
G-RAG: Knowledge Expansion in Material Science

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In the field of Material Science, effective information retrieval systems are essential
2 for facilitating research. Traditional Retrieval-Augmented Generation (RAG)
3 approaches in Large Language Models (LLMs) often encounter challenges such
4 as outdated information, hallucinations, limited interpretability due to context
5 constraints, and inaccurate retrieval. To address these issues, Graph RAG integrates
6 graph databases to enhance the retrieval process. Our proposed method processes
7 Material Science documents by extracting key entities (referred to as MatIDs)
8 from sentences, which are then utilized to query external Wikipedia knowledge
9 bases (KBs) for additional relevant information. We implement an agent-based
10 parsing technique to achieve a more detailed representation of the documents. Our
11 improved version of Graph RAG called G-RAG further leverages a graph database
12 to capture relationships between these entities, improving both retrieval accuracy
13 and contextual understanding. This enhanced approach demonstrates significant
14 improvements in performance for domains that require precise information retrieval,
15 such as Material Science.

16 1 Introduction

17 LLMs exhibit impressive capabilities but encounter challenges such as hallucinations, outdated
18 information, and untraceable, opaque reasoning. The RAG approach addresses these issues by
19 combining the strengths of LLMs with the vast, continuously updated resources of external databases
20 [1]. Graph-enhanced RAG methods build on this by leveraging rich semantic interconnections and
21 relational data, enabling more precise entity linking, enhanced semantic context, and improved
22 knowledge extraction [2, 3]. Additionally, researchers have introduced innovative graph-based
23 context adaptation techniques that refine word embeddings to better capture semantic relationships,
24 consistently outperforming traditional methods in various Natural Language Processing (NLP) tasks
25 [4, 5]. Graph-based RAG provides a more nuanced and accurate representation of complex domains,
26 enabling LLMs to generate responses with enhanced factual precision and contextual relevance [6, 7].
27 This capability is especially valuable for domain-specific applications in fields such as material
28 science and biomedicine, where accurate and detailed information is crucial [8, 9, 10]. Serving as
29 a domain-specific knowledge server, the Semantic Context Enhancer extracts and delivers detailed
30 descriptions of relevant concepts and entities, including their interrelationships, thereby equipping
31 the LLM with a deeper semantic understanding [10]. Additionally, leveraging graph structures to
32 improve knowledge retrieval and response generation, as exemplified by methods like AriGraph,
33 has shown significant enhancements in decision-making and planning capabilities [11]. This study
34 explores the improvement of information retrieval and knowledge generation in complex, specialized
35 domains through the integration of the G-RAG pipeline, addressing limitations of existing approaches
36 and advancing performance in targeted fields.

37 2 Methodology

38 The retrieval process of Naive RAG includes a diverse range of MatIDs, which ensures variety but can
39 also introduce less relevant information. This issue can be mitigated through prompt engineering in
40 the RAG configuration, allowing the LLM to continue generating accurate responses [12]. However,
41 there are two main limitations to this approach. First, LLMs have a fixed context window, which
42 restricts the number of tokens they can process simultaneously. This limitation hinders the model’s
43 ability to manage large volumes of retrieved data effectively [13], especially when the dataset is
44 extensive and varied. Despite advancements like Google’s Gemini, which uses a caching system to
45 handle extended contexts, the fixed context window of LLMs remains a significant constraint [14, 15].
46 Although providing the model with more relevant information might seem beneficial, increasing
47 the context length does not necessarily improve the accuracy of information retrieval or response
48 generation [16]. This problem becomes even more pronounced when the retrieved context includes a
49 mix of diverse but only marginally relevant data, potentially diluting the focus on the critical entities
50 or concepts needed for an accurate response [17]. This is where Graph RAG proves to be valuable, as
51 it enhances the retrieval process by focusing on the most relevant information.

52 2.1 Graph RAG vs G-RAG

53 Graph RAG effectively merges the strengths of retrieval-based and generative methods to enhance
54 LLMs’ capability to generate accurate, relevant, and contextually enriched responses [18]. While
55 supplying an LLM with text chunks from extensive documents may result in issues with context,
56 factual precision, and language coherence, Graph RAG addresses these limitations by utilizing a
57 knowledge graph as a source of structured, factual information [19]. The knowledge graph provides
58 detailed entity information, including attributes and relationships, allowing the LLM to gain a deeper
59 understanding and produce more informed, precise responses. In our G-RAG system, entity linking
60 is a fundamental component, enabling the extraction of specific entities (key terms or concepts) from
61 the text using an entity extractor like a Span Parser. These identified entities are then used to query
62 an external retriever, which fetches relevant MatIDs and their corresponding information from a
63 Wikipedia knowledge base [20]. This targeted retrieval process ensures that the selected MatIDs are
64 highly relevant and accurate, thereby preserving the integrity of the constructed knowledge graph [21].
65 Following this, an LLM formulates a query that is sent to the graph database. The graph database
66 retrieves relevant information, which is processed by the LLM to generate a final, comprehensive
67 response. The complete architecture of our G-RAG system is illustrated in Figure 1.

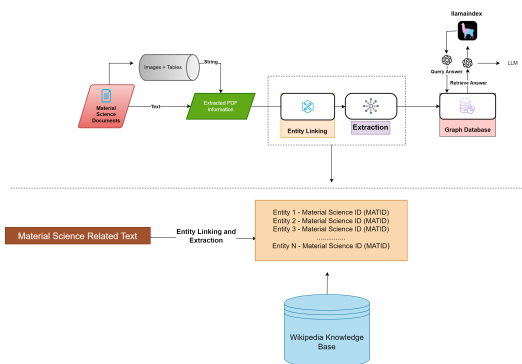


Figure 1: Architecture of G-RAG System

68 2.2 PDF Parsing

69 We parse PDFs by categorizing their content into text, figures, and tables. For figure extraction, we
70 employ the Phi-3.5 Vision Instruct model, specifically tailored to identify material science-related
71 images using a vision agent system. We utilize Microsoft’s Table Transformer in the tabular data
72 extraction process. Furthermore, we apply a smart chunking technique to enhance the precision of
73 data segmentation [22]. Accurate parsing is essential for subsequent tasks such as Entity Linking,
74 Relation Extraction, and Graph Retrieval Augmented Generation, as it ensures the accuracy and

75 relevance of the answers retrieved from the database. Appendix A.1 provides a detailed overview of
76 our document parsing process.

77 2.3 Entity Linking and Relation Extraction

78 Entity Linking (EL) refers to the process of mapping ambiguous mentions in a text to specific,
79 identifiable named entities within a knowledge base [23]. It involves recognizing all potential entities
80 mentioned in the given input and accurately associating them with corresponding entries in a reference
81 knowledge base, such as Wikipedia. Relation Extraction (RE) refers to the process of identifying and
82 classifying semantic relationships between entities mentioned within a given text. This task involves
83 mapping the detected entities to specific relation categories defined in a reference knowledge base,
84 such as Wikipedia. The entity linking and relation extraction process is depicted in Appendix A.2.

85 2.4 Span Parser

86 The Span Parser module functions as our G-RAG system’s initial information retrieval component,
87 employing an approach inspired by the Retrieval Process [24]. This module operates on the principle
88 of semantic similarity between the current knowledge base (KB) and a comprehensive collection
89 of textual passages (Wikipedia Database) representing entities and relations. At its core, the Span
90 Parsing module utilizes an encoder to generate dense vector representations of both the knowledge
91 base (KB) q and each passage p in the additional knowledge base collection. These representations,
92 denoted as $E(q)$ and $E(p)$ respectively, are high-dimensional embeddings that capture the semantic
93 content of the text. The module computes a similarity score between the current Knowledge Base
94 and additional Knowledge Base (Wikipedia data) using a dot product operation, yielding the most
95 relevant relations with respect to the extended knowledge base q :

$$\text{sim}(q, p) = E(q)^\top \cdot E(p)$$

96 This score quantifies the relevance of each passage of the additional knowledge base to the given cur-
97 rent KB passage’s sentence, enabling the module to rank and retrieve the most pertinent information.

98 2.5 Passage Processor

99 The Passage Processor (PP) component in our G-RAG system employs a unified approach to process
100 the existing knowledge base and retrieved passages. Given a current Knowledge Base (KB) Q and
101 a set of N retrieved passages $\{P_1, \dots, P_n\}$, the Passage Processor constructs chunks of current
102 KB. In each chunk, we utilize each input sequence $S = [Q; \tau_0; P_1; \tau_1; \dots; P_n; \tau_n]$, where τ_i are
103 delimiter tokens. This sequence is encoded using a Transformer model T , producing contextual
104 embedding $E = T(S)$. The Passage Processor subsequently identifies relevant spans within Q
105 through a two-stage process [24]. Initially, it computes start probabilities $P^s(q_i)$ for each token q_i in
106 Q using a learned function $f^s(E)$. Subsequently, for each potential start position s , it calculates end
107 probabilities $P^e(q_j | s)$ for tokens q_j (where $j \geq s$) using another learned function $f^e(E, s)$. This
108 formulation enables the prediction of overlapping spans, enhancing the model’s capability to handle
109 complex queries. During the process, spans (s, e) are predicted if $P^s(q_s) > \theta_s$ and $P^e(q_e | s) > \theta_e$,
110 where θ_s and θ_e are predefined thresholds. This design enables the Passage Processor to process the
111 entire knowledge base chunk by chunk efficiently, identifying relevant text spans for downstream
112 tasks such as entity linking and relation extraction.

113 3 Experimental Settings

114 Our dataset consists of ten carefully designed handwritten queries, aimed at evaluating and differen-
115 tiating the capabilities of various RAG systems. Sample queries from this dataset are presented in
116 Appendix A.3. To evaluate the performance of RAG systems, we employ various metrics, including
117 correctness, faithfulness, context, and answer relevancy scores. Correctness assesses the accuracy
118 of the generated response, while faithfulness evaluates the factual accuracy based on the retrieved
119 documents. Finally, the context and answer relevancy score measures how well the response aligns
120 with the given query. A detailed description of these evaluation metrics is provided in Appendix
121 A.4. For entity linking and relation extraction, we use the relik-entity-linking-large model [25],
122 while the jina-embeddings-v2-base-en model [26], with a sequence length of 8192, is employed for

123 embeddings. Additionally, we utilize LLama 3.1 8B and LLama 3.1 70B as large language models,
124 both of which produce comparable results.

125 4 Results and Discussion

126 This section presents all of our experimental results. We conducted the computational tasks using
127 the NVIDIA Tesla A100 Ampere 40 GB GPU. The performance of the Naive RAG, Graph RAG,
128 and the G-RAG system was evaluated using our dataset. Appendix A.5 provides example queries
129 and the corresponding responses from the RAG systems, evaluated across different metrics. The
130 experimental results are summarized in Table 1.

Table 1: Experimental Results

Pipeline	Score	No. of queries	Mean	Standard Deviation
Naive RAG	Correctness	10	2.43	1.51
	Faithfulness		0.70	0.48
	Relevancy		0.39	0.28
Graph RAG	Correctness	10	3.30	2.00
	Faithfulness		0.90	0.32
	Relevancy		0.18	0.26
G-RAG	Correctness	10	3.90	1.10
	Faithfulness		0.90	0.32
	Relevancy		0.34	0.32

131 The comparative analysis of three RAG pipelines - Vector/Naive RAG, G-RAG, and Graph RAG -
132 showed interesting patterns in their performance across three critical dimensions. A one-way Analysis
133 of Variance (ANOVA) as described in Appendix A.4.5 was performed, examining correctness
134 $F(2, 24) = 2.39, p = 0.113$, faithfulness $F(2, 27) = 1.04, p = 0.368$, and context and answer
135 relevancy $F(2, 27) = 1.04, p = 0.368$. While no statistically significant differences were found at the
136 standard significance level ($\alpha = 0.05$), the descriptive statistics highlighted meaningful variations in
137 performance. Specifically, Vector/Naive RAG outperformed the others in terms of context relevancy,
138 with a mean score of 0.3875. This was followed by G-RAG (mean score of 0.3375), while Graph
139 RAG exhibited the lowest mean score of 0.1750. The substantial standard deviations observed across
140 all pipelines, ranging from 0.2630 to 0.3162, suggest notable performance variability depending on
141 the query. This variability highlights the challenge of consistency in RAG systems. The superior
142 performance of G-RAG over the basic Graph RAG can be attributed to the inclusion of a material
143 science knowledge base, emphasizing the critical role of domain-specific knowledge in enhancing
144 model accuracy. The superior context relevancy performance of the traditional Vector/Naive RAG
145 challenges the assumption that graph-based approaches inherently provide better retrieval capabilities.
146 G-RAG has proven to be a well-rounded solution, effectively balancing the metrics of correctness,
147 relevancy, and faithfulness. The significant drop in relevancy scores for Graph RAG highlights the
148 critical role of entity linking in G-RAG’s design. This suggests that the effectiveness of knowledge
149 integration mechanisms, including entity linking, plays a substantial role in improving retrieval
150 performance. These findings indicate that while graph-based approaches show promise, their success
151 heavily depends on the quality of knowledge integration and the sophistication of the entity-linking.

152 5 Conclusion and Future Work

153 Our findings indicate that integrating graph-based techniques and ensuring robust entity linking
154 with external databases can significantly enhance the performance of the Graph RAG pipeline,
155 particularly in terms of response relevance and accuracy. This approach also mitigates the challenge
156 of maintaining relevance observed in standard Graph RAG implementations. Future work could
157 include developing a larger knowledge base tailored to material science as an extended information
158 source, as well as creating a material science-specific entity linking model. Additionally, establishing
159 a comprehensive evaluation metric for Graph RAG would provide deeper insights into the process
160 and its effectiveness.

References

- [1] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey, 2024.
- [2] Sanat Sharma, David Seunghyun Yoon, Franck Deroncourt, Dewang Sultania, Karishma Bagga, Mengjiao Zhang, Trung Bui, and Varun Kotte. Retrieval augmented generation for domain-specific question answering. *arXiv preprint arXiv:2404.14760*, 2024.
- [3] Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, H. Qu, and Jian Guo. Think-on-graph 2.0: Deep and interpretable large language model reasoning with knowledge graph-guided retrieval. *arXiv preprint arXiv:2407.10805*, 2024.
- [4] Tanvi Sandhu and Ziad Kobti. Exploration of word embeddings with graph-based context adaptation for enhanced word vectors. In *Proceedings of the International Florida Artificial Intelligence Research Society Conference*, 2024.
- [5] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization, 2024.
- [6] What Is Graph RAG? — ontotext.com. <https://www.ontotext.com/knowledgehub/fundamentals/what-is-graph-rag/>. [Accessed 02-09-2024].
- [7] Chansol Park, Hayoung Lee, and Ok-Ran Jeong. Leveraging medical knowledge graphs and large language models for enhanced mental disorder information extraction. *Future Internet*, 2024.
- [8] Markus J. Buehler. Generative retrieval-augmented ontologic graph and multiagent strategies for interpretive large language model-based materials design. *ACS Engineering Au*, 2024.
- [9] Julien Delile, Srayanta Mukherjee, A. V. Pamel, and Leonid Zhukov. Graph-based retriever captures the long tail of biomedical knowledge. *arXiv preprint arXiv:2402.12352*, 2024.
- [10] Ciyuan Peng, Feng Xia, Mehdi Naseriparsa, and Francesco Osborne. Knowledge graphs: Opportunities and challenges, 2023.
- [11] Petr Anokhin, Nikita Semenov, A. N. Sorokin, Dmitry Evseev, Mikhail Burtsev, and Evgeny Burnaev. Arigraph: Learning knowledge graph world models with episodic memory for llm agents. *arXiv preprint arXiv:2407.04363*, 2024.
- [12] Thomas Merth, Qichen Fu, Mohammad Rastegari, and Mahyar Najibi. Superposition prompting: Improving and accelerating retrieval-augmented generation, 2024.
- [13] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of large language models via positional interpolation, 2023.
- [14] Cunchen Hu, Heyang Huang, Junhao Hu, Jiang Xu, Xusheng Chen, Tao Xie, Chenxi Wang, Sa Wang, Yungang Bao, Ninghui Sun, and Yizhou Shan. Memserve: Context caching for disaggregated llm serving with elastic memory pool, 2024.
- [15] Xiaohua Wang, Zhenhua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Li Wang, Shizheng Li, Qian Qi, et al. Searching for best practices in retrieval-augmented generation. *arXiv preprint arXiv:2407.01219*, 2024.
- [16] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts, 2023.
- [17] Anand Subramanian. Building a Biomedical Entity Linker with LLMs — towardsdatascience.com. <https://towardsdatascience.com/building-a-biomedical-entity-linker-with-llms-d385cb85c15a>. [Accessed 02-09-2024].

- 207 [18] Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang,
208 and Siliang Tang. Graph retrieval-augmented generation: A survey. *arXiv preprint*
209 *arXiv:2408.08921*, 2024.
- 210 [19] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven
211 Truitt, and Jonathan Larson. From local to global: A graph rag approach to query-focused
212 summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- 213 [20] Wikipedia Knowledge Base. [relik-ie/relik-reader-deberta-v3-large-re-wikipedia](https://www.wikidata.org/wiki/Relik-IE/Relik-Reader-DeBERTa-v3-Large-re-Wikipedia).
214 [Accessed 04-09-2024].
- 215 [21] Entity Linking and Relationship Extraction With Relik in Lla-
216 maIndex — neo4j.com. [https://neo4j.com/developer-blog/
217 entity-linking-relationship-extraction-relik-llamaindex/](https://neo4j.com/developer-blog/entity-linking-relationship-extraction-relik-llamaindex/). [Accessed
218 02-09-2024].
- 219 [22] Darren Oberst, MacOS, Jeff Turnham, Jessica Berliner, Will Taner, Prashant Rajesh Iyer,
220 NYDocutest, Adelina Jiang, Aryan Chauhan, Christopher Harrison, SNEHA KUMARI, Viren-
221 dra Singh, simonrisman, shneeba, Vedant Sudhir Patil, Rohan Sharma, Raghav Dixit, Rahul
222 Khandait, Uğur Çekmez, Ayushri, adithyasudhan, Chelsea, Ujjwal Jha, momodingaling,
223 osi1880vr, Kamalakar Satapathi, philipkd, Rajesh Adhikari, Saurav Kumar Mahato, and Seva
224 Skvortsov. *llmware-ai/llmware*, 9 2024.
- 225 [23] Simone Tedeschi, Simone Conia, Francesco Cecconi, and Roberto Navigli. Named entity
226 recognition for entity linking: What works and what’s next. In Marie-Francine Moens, Xuanjing
227 Huang, Lucia Specia, and Scott Wen tau Yih, editors, *EMNLP (Findings)*, pages 2584–2596.
228 Association for Computational Linguistics, 2021.
- 229 [24] Riccardo Orlando, Pere-Lluis Huguet-Cabot, Edoardo Barba, and Roberto Navigli. Relik:
230 Retrieve and link, fast and accurate entity linking and relation extraction on an academic budget,
231 2024.
- 232 [25] Relik Entity Linking Model. [https://huggingface.co/sapienzanlp/
233 relik-entity-linking-large](https://huggingface.co/sapienzanlp/relik-entity-linking-large). [Accessed 02-09-2024].
- 234 [26] Jina Embeddings Model. [https://huggingface.co/jinaai/
235 jina-embeddings-v2-base-en](https://huggingface.co/jinaai/jina-embeddings-v2-base-en). [Accessed 02-09-2024].

236 **A Appendix**

237 **A.1 Documents Parsing Method**

238 This section illustrates our document parsing pipeline, as shown in Figures 2 and 3. Efficient document
 239 parsing is crucial for enabling RAG systems to generate responses with high factual accuracy and
 240 precision.

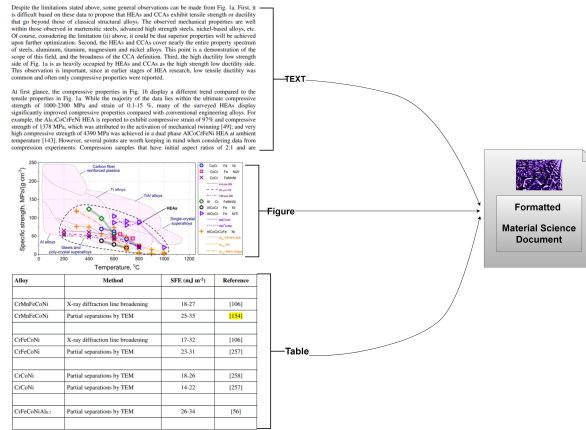


Figure 2: Document Parsing

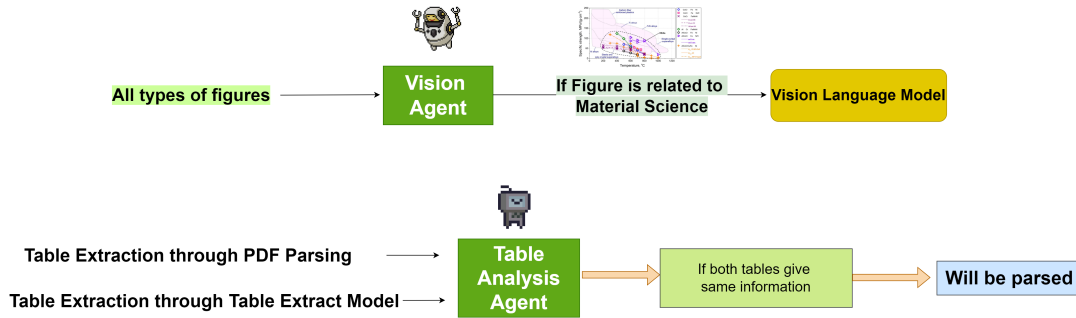


Figure 3: Validity Check by Agent System

241 **A.2 Entity Linking and Relation Extraction**

242 In this section, we provide a visual representation of the entity linking and relation extraction process,
 243 as depicted in Figures 4, 5, 6, and 7. These processes are essential components of our G-RAG system.

244 **Coreference Resolution:** Coreference resolution, mentioned in Figure 4 involves identifying different
 245 expressions in a text that refer to the same entity. This process is crucial for understanding the
 246 relationships between various mentions of an entity within a given context.

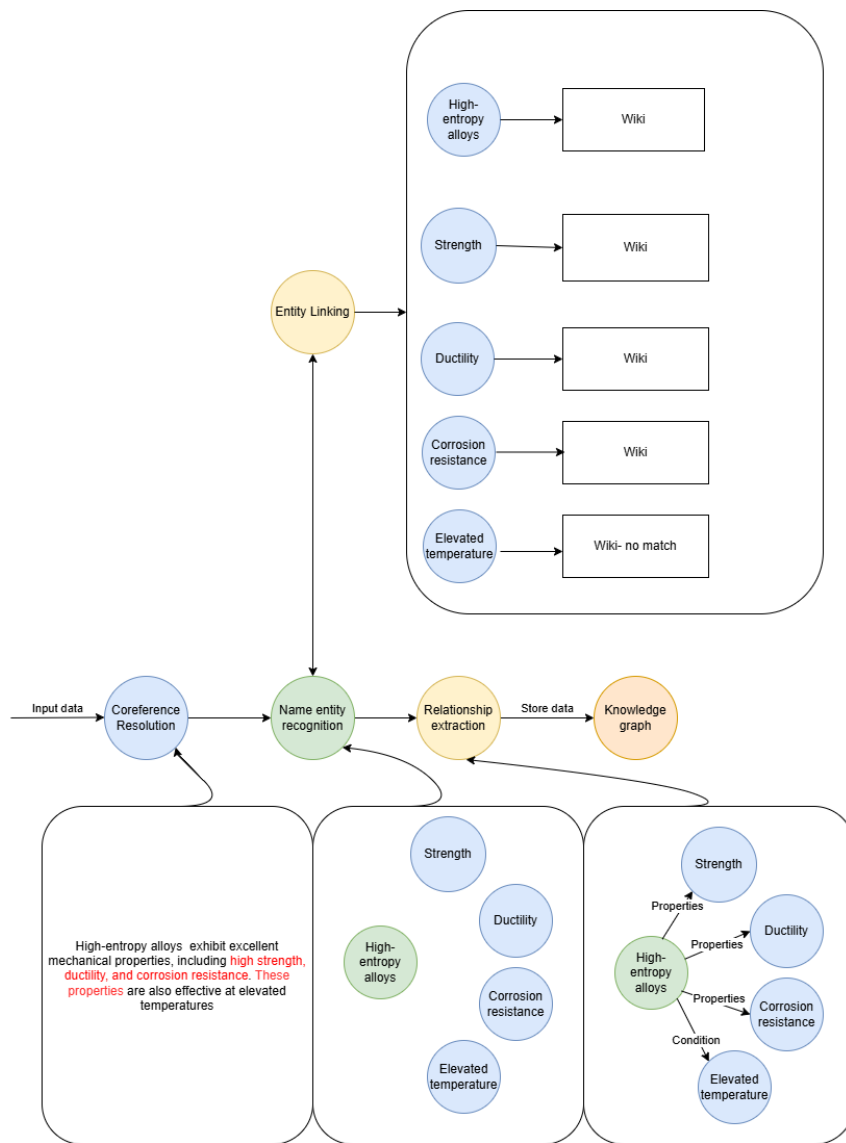


Figure 4: Entity Linking and Relation Extraction

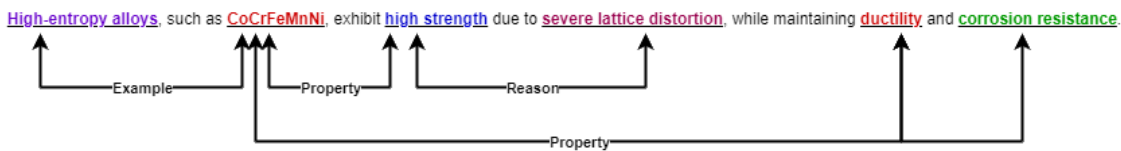


Figure 5: Entity Linking

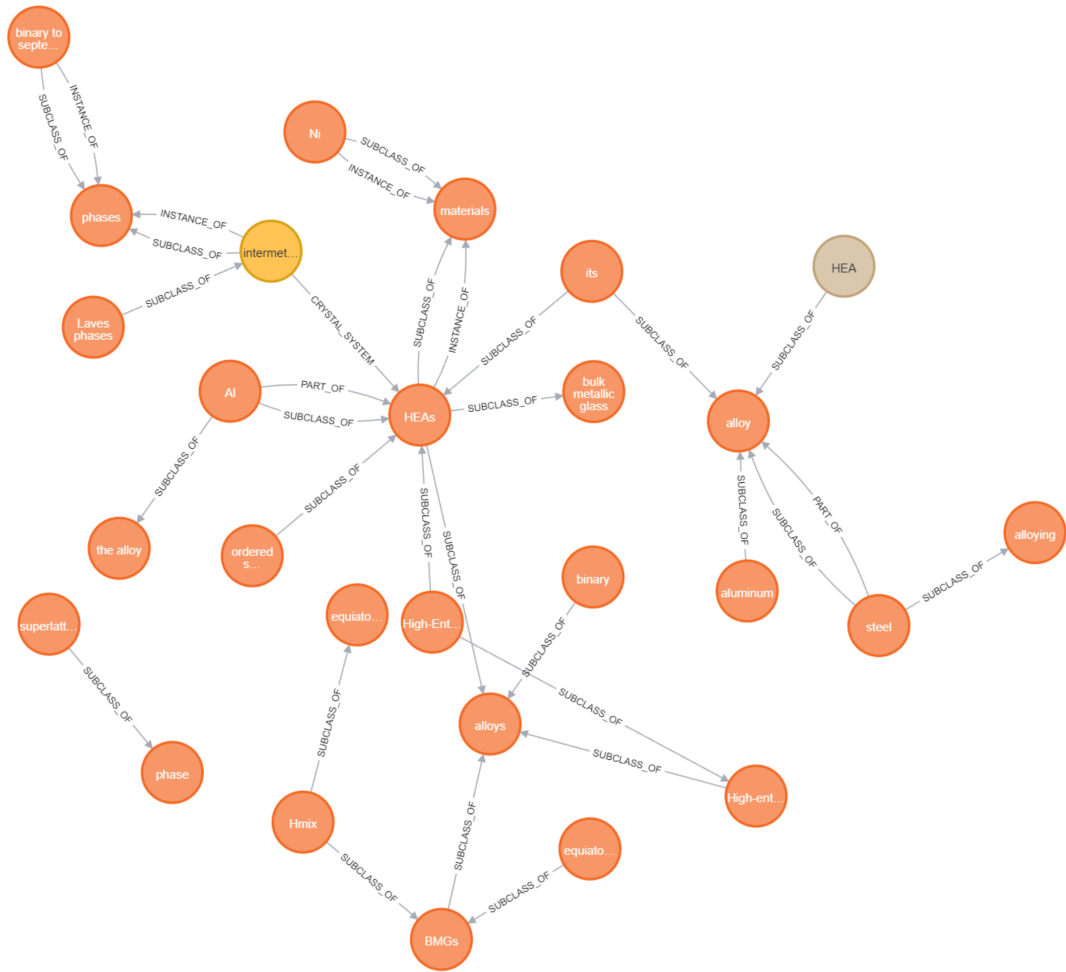


Figure 6: Relationship among Various High-entropy alloy Components

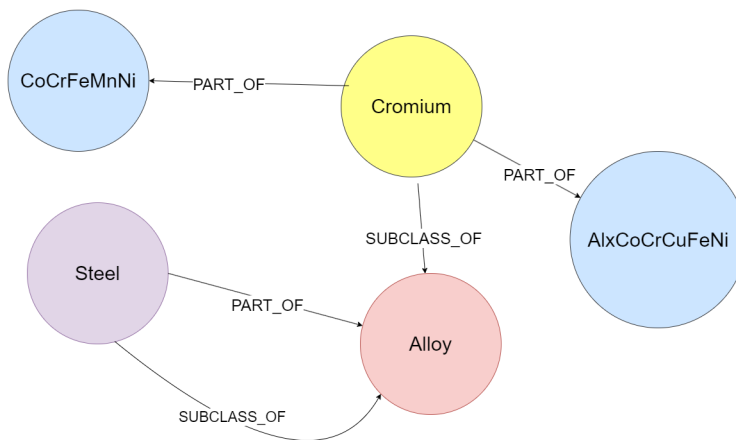


Figure 7: Another Relationship among Various High-entropy alloy Components

247 **A.3 Examples from Our Dataset**

248 In this section, we present sample queries from our dataset in Table 2, covering a range from simple
 249 to more complex queries.

Table 2: Example Queries

Query	What is the yield strength of the CrMnFeCoNi alloy at 600 K, 700 K with 4 μm grain size?
Ground Truth	The yield strength of the CrMnFeCoNi alloy at 600 K is 290 MPa, and at 700 K, it is 285 MPa.
Query	What is the CRSS of CrMnFeCoNi at the tension in room temperature?
Ground Truth	The Critical Resolved Shear Stress (CRSS) of the CrMnFeCoNi alloy has been measured at 53 MPa at room temperature and 175 MPa at 77 K.
Query	What is the stacking fault energy of CrCoNi?
Ground Truth	The stacking fault energy of CrCoNi is $18 - 26 \text{ mJ/m}^2$.
Query	At room temperature, what is the Hall-Petch slope of the cantor alloy?
Ground Truth	At room temperature, the Hall-Petch slope of the cantor alloy was determined to be $494 \text{ MPa } \mu\text{m}^{-1/2}$.
Query	What is the stacking fault energy of the cantor alloy?
Ground Truth	The stacking fault energy of the cantor alloy was estimated to be $\sim 30 \text{ mJ m}^{-2}$.
Query	What is the yield strength and ultimate tensile strength of TiZrNbHfTa after 1000°C annealing?
Ground Truth	After 1000 °C, the yield strength will be 1145 MPa, and the ultimate tensile strength will be 1262 MPa.
Query	What is the CRSS of CrFeCoNiAl0.3 in compression at room temperature?
Ground Truth	CRSS of CrFeCoNiAl0.3 in compression at room temperature is 54 MPa.

250 **A.4 LLM RAG Evaluation Metrics**

251 This section provides detailed descriptions of the various evaluation metrics used for RAG systems.

252 **A.4.1 Correctness**

253 Given a query q , a generated answer g , and an optional reference answer r , the
 254 `CorrectnessEvaluator` computes a score s using an LLM. This score is then compared against a
 255 threshold T to determine whether the generated answer is correct or passing.

Prompt	Constructed from q, g, r
$E(g, q, r)$	LLM Response to Prompt
$(s, reasoning)$	<code>parser_function(E(g, q, r))</code>
passing	$s \geq T$
EvaluationResult	$\{q, g, passing, s, reasoning\}$

256 **A.4.2 Faithfulness Evaluation**

257 Given a query q , a response r , and a set of context documents C , the `FaithfulnessEvaluator`
 258 performs the following steps:

Context Documents	Transform C into Document objects
Index	Create <code>SummaryIndex</code> from Document objects
Query Engine	Create query engine using LLM, <code>eval_template</code> , and <code>refine_template</code>
Evaluation	Perform a query on the response using the query engine
Raw Response	Obtain <code>raw_response_txt</code> from the query engine
Passing	$\begin{cases} \text{True} & \text{if yes is found in raw_response_txt} \\ \text{False} & \text{otherwise} \end{cases}$
Score	$\begin{cases} 1.0 & \text{if passing is True} \\ 0.0 & \text{otherwise} \end{cases}$
Feedback	<code>raw_response_txt</code>

259 The evaluation result is given by:

$$\text{EvaluationResult} = \{q, r, C, passing, score, feedback\}$$

260 **A.4.3 Answer Relevancy**

261 Let q be the query, r the response, and $\{c_1, c_2, \dots, c_n\}$ the contexts. Define the following:

$$\text{Documents} = \{d_i \mid d_i = \text{Document}(\text{text} = c_i) \text{ for } i = 1, 2, \dots, n\}$$

262

$$\text{Index} = \text{SummaryIndex}(\text{Documents})$$

263

$$\text{query_response} = \text{Question: } q \text{ Response: } r$$

264 Evaluate the query-response pair with:

$$\text{response_obj} = \text{QueryEngine}(\text{Index}).\text{aquery}(\text{query_response})$$

265 Let:

$$\text{raw_response_txt} = \text{str}(\text{response_obj})$$

266 Then:

$$\text{passing} = \begin{cases} \text{True} & \text{if "yes" is in raw_response_txt.lower()} \\ \text{False} & \text{otherwise} \end{cases}$$

267

$$\text{score} = \begin{cases} 1.0 & \text{if passing} \\ 0.0 & \text{otherwise} \end{cases}$$

268 *The output is:*

EvaluationResult = {*q*, *r*, passing, score, feedback = raw_response_txt, contexts = {*c*₁, . . . , *c*_{*n*}}

269 **A.4.4 Context Relevancy**

270 *Let q* be the query, {*c*₁, *c*₂, . . . , *c*_{*n*}} be the contexts. Define:

Documents = {*d*_{*i*} | *d*_{*i*} = Document(text = *c*_{*i*})}

271

Index = SummaryIndex(Documents)

272 *Evaluate the query q* using:

query_engine = Index.as_query_engine(llm, eval_template, refine_template)

273

response_obj = query_engine.aquery(*q*)

274 *Let:*

raw_response_txt = str(response_obj)

275 *Parse the result:*

score, reasoning = parser_function(raw_response_txt)

276 *Score threshold:*

score_threshold = 4.0

277 *Calculate:*

$$\text{score} = \frac{\text{score}}{\text{score_threshold}}$$

278 *Return:*

EvaluationResult = {*q*, {*c*₁, . . . , *c*_{*n*}}, score, feedback = raw_response_txt, invalid_result, invalid_reason}

279 **A.4.5 Analysis of Variance (ANOVA)**

280 ANOVA is a fundamental statistical method used to compare means across multiple groups to
281 determine if there are statistically significant differences between them. This study utilizes a one-
282 way ANOVA, which examines the effect of a single independent variable - in this case, the type
283 of RAG pipeline - on a dependent variable (performance metrics). The mean score reflects the
284 average performance of each method across all 10 queries, offering an overall assessment of its
285 effectiveness for the given metrics. A mean score closer to the highest possible value suggests that the
286 method consistently delivers superior results, indicating strong performance across various queries.
287 Conversely, a lower mean score points to weaker overall performance, highlighting areas where the
288 method may be less effective. Essentially, the mean score serves as a summary indicator of each
289 method's typical efficacy, providing a clear comparison of their relative strengths in achieving the
290 desired outcomes.

291 The F-statistic in ANOVA quantifies the ratio of variance between groups to variance within groups,
292 with larger F-values indicating greater differences among the groups. The degrees of freedom (df) are
293 denoted as F(2, 24) for correctness and F(2, 27) for faithfulness and relevancy, indicating the number
294 of independent values that can vary in the calculation. Here, the first value (2) represents the degrees
295 of freedom for the groups (number of groups minus 1), and the second value (24/27) represents the
296 degrees of freedom for the error term (total observations minus the number of groups). The p-value
297 assesses the probability of observing such differences between groups. Typically, a p-value below the
298 significance level ($\alpha = 0.05$) suggests statistically significant differences between the groups.

299 **A.5 Examples of Different Metrics**

300 This section presents three example queries along with their corresponding responses from the RAG
 301 systems. Providing optimal responses to these queries requires effective information retrieval from
 302 text, figures, and tables. Additionally, we illustrate the application of various RAG evaluation metrics
 303 through these examples. Table 3 compares the performance using high-entropy alloy documents,
 304 demonstrating that G-RAG delivers contextually more relevant responses. The query in Table 4
 305 necessitates accurate information extraction from Figure 8. The responses show that G-RAG with the
 306 parsing method outperforms other RAG systems, providing results that are closely aligned with the
 307 ground truth values. Finally, the query in Table 5 requires accurate value extraction from the table
 308 presented in Figure 9. In this case, G-RAG with the parsing method also outperforms other RAG
 309 systems. These results demonstrate that the integration of an efficient parsing method significantly
 310 enhances the performance of the G-RAG system.

Table 3: Comparison using **High-entropy alloy Documents (Context Relevancy)**

Query	What is Chromium?
G-RAG	Chromium is often a significant component in High-Entropy alloys (HEAs), particularly in the CoCrFeNi alloy, where it enhances corrosion resistance in both NaCl and H ₂ SO ₄ solutions. However, when Chromium is replaced with Manganese in the Co-Cr-Fe-Ni alloy, the corrosion resistance in H ₂ SO ₄ solution decreases, becoming inferior to 304 SS. The addition of certain elements like Al, Cu, B, and Mo can alter the corrosion resistance of various HEAs, often due to the formation of new phases or the impact on the existing phases' composition. The relationship between the degree of substitution (δ) and the enthalpy of mixing (H_{mix}) for intermetallic compounds is an area of interest, potentially implying that as composition increases, the enthalpy of mixing decreases, which could have implications for the thermodynamics of alloys and intermetallic compound formation.
Naive RAG	Chromium is a chemical element with the symbol Cr and atomic number 24. It is a hard, silver-white, and highly corrosion-resistant metal that is used in a variety of applications, including stainless steel, chrome plating, and catalytic converters in vehicles.
Graph RAG	Chromium is a chemical element with the symbol Cr and atomic number 24. It is a steely-grey, lustrous, hard, and brittle transition metal. Chromium is also the main additive in stainless steel, to which it adds anti-corrosive properties. It is highly valued as a metal that is able to be highly polished while resisting tarnishing. Chromium is also used in various industrial applications, including the production of stainless steel, and chrome plating, and as a corrosion-resistant coating.

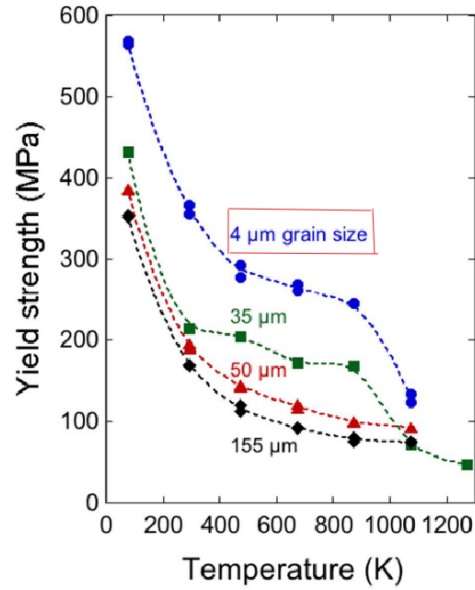


Figure 8: Information Extraction from Graph Images

Table 4: Comparison of Correct Information Extraction from Graph

Query	What is the yield strength of the CrMnFeCoNi alloy at 600 K, 700 K with 4 μm grain size?
Ground Truth	The yield strength of the CrMnFeCoNi alloy at 600 K is 290 MPa, and at 700 K, it is 285 MPa.
(G-RAG + Parsing Method)	The yield strength of CrMnFeCoNi alloy at 600 K is 310 MPa, and at 700 K is around 300 MPa.
Naive RAG	Unfortunately, the information provided does not contain any data about the yield strength of the CrMnFeCoNi alloy at 600 K, 700 K, or with a 4 μm grain size.
G-RAG	The yield strength of CrMnFeCoNi alloy at 600 K is 450 MPa, and at 700 K is 400 MPa.

Alloy	Loading	Temperature	CRSS (MPa)	Reference
CrMnFeCoNi	Tension	RT	53 - 60	[228]
CrMnFeCoNi	Tension	RT	63	[229]
CrMnFeCoNi	Tension	RT	70	[224]
CrMnFeCoNi	Tension	RT	78 - 82	[225]
CrMnFeCoNi	Tension	RT	82	[230]
CrMnFeCoNi	Compression ¹	RT	33 - 43	[154]
CrMnFeCoNi	Tension	77 K	13 5	[229]
CrMnFeCoNi	Tension	77 K	15 5	[230]
CrMnFeCoNi	Tension	77 K	14 5- 17 2	[228]
CrMnFeCoNi	Tension	77 K	17 5	[224]
CrMnFeCoNi	Tension	203 K	10 3	[230]
CrMnFeCoNi	Tension	243 K	92	[230]

Figure 9: Information Extraction from Table Images

Table 5: Comparison of Correct Information Extraction from Table

Query	What is the CRSS of CrMnFeCoNi at the tension in room temperature?
Ground Truth	The Critical Resolved Shear Stress (CRSS) of the CrMnFeCoNi alloy has been measured at 53 MPa at room temperature and 175 MPa at 77 K.
(G-RAG + Parsing Method)	The CRSS of the CrMnFeCoNi alloy has been determined to be 53 MPa at room temperature and 175 MPa at 77 K.
Naive RAG	Unfortunately, the provided context does not contain information about the CRSS (Critical Resolved Shear Stress) of CrMnFeCoNi at tension at room temperature.
G-RAG	Sorry, I do not have any information about the CRSS of CrMnFeCoNi at the tension at room temperature.

311 **NeurIPS Paper Checklist**

312 **1. Claims**

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

315 Answer: [\[Yes\]](#)

316 Justification: Our main claims made in the abstract and introduction accurately reflect the
317 paper's contribution and scope.

318 Guidelines:

- 319 • The answer NA means that the abstract and introduction do not include the claims
320 made in the paper.
- 321 • The abstract and/or introduction should clearly state the claims made, including the
322 contributions made in the paper and important assumptions and limitations. A No or
323 NA answer to this question will not be perceived well by the reviewers.
- 324 • The claims made should match theoretical and experimental results, and reflect how
325 much the results can be expected to generalize to other settings.
- 326 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
327 are not attained by the paper.

328 **2. Limitations**

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

330 Answer: [\[Yes\]](#)

331 Justification: We have discussed the limitations of the work as well as future improvements
332 that can be made to our work.

333 Guidelines:

- 334 • The answer NA means that the paper has no limitation while the answer No means that
335 the paper has limitations, but those are not discussed in the paper.
- 336 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 337 • The paper should point out any strong assumptions and how robust the results are to
338 violations of these assumptions (e.g., independence assumptions, noiseless settings,
339 model well-specification, asymptotic approximations only holding locally). The authors
340 should reflect on how these assumptions might be violated in practice and what the
341 implications would be.
- 342 • The authors should reflect on the scope of the claims made, e.g., if the approach was
343 only tested on a few datasets or with a few runs. In general, empirical results often
344 depend on implicit assumptions, which should be articulated.
- 345 • The authors should reflect on the factors that influence the performance of the approach.
346 For example, a facial recognition algorithm may perform poorly when image resolution
347 is low or images are taken in low lighting. Or a speech-to-text system might not be
348 used reliably to provide closed captions for online lectures because it fails to handle
349 technical jargon.
- 350 • The authors should discuss the computational efficiency of the proposed algorithms
351 and how they scale with dataset size.
- 352 • If applicable, the authors should discuss possible limitations of their approach to
353 address problems of privacy and fairness.
- 354 • While the authors might fear that complete honesty about limitations might be used by
355 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
356 limitations that aren't acknowledged in the paper. The authors should use their best
357 judgment and recognize that individual actions in favor of transparency play an impor-
358 tant role in developing norms that preserve the integrity of the community. Reviewers
359 will be specifically instructed to not penalize honesty concerning limitations.

360 **3. Theory Assumptions and Proofs**

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

363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415

Answer: [Yes]

Justification: We have provided the full set of assumptions and a complete proof.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

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

Answer: [Yes]

Justification: The authors of this paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

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

419 Answer: [Yes]

420 Justification: We have a proper codebase of everything we do to make it as reproducible as
421 possible. But for reproducing it, if anyone wants to create a new knowledge base with the
422 given documents, it is mandatory to have at least 24 GB VRAM and 60 GB System RAM.

423 Guidelines:

- 424 • The answer NA means that paper does not include experiments requiring code.
- 425 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
426 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 427 • While we encourage the release of code and data, we understand that this might not be
428 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
429 including code, unless this is central to the contribution (e.g., for a new open-source
430 benchmark).
- 431 • The instructions should contain the exact command and environment needed to run to
432 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
433 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 434 • The authors should provide instructions on data access and preparation, including how
435 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 436 • The authors should provide scripts to reproduce all experimental results for the new
437 proposed method and baselines. If only a subset of experiments are reproducible, they
438 should state which ones are omitted from the script and why.
- 439 • At submission time, to preserve anonymity, the authors should release anonymized
440 versions (if applicable).
- 441 • Providing as much information as possible in supplemental material (appended to the
442 paper) is recommended, but including URLs to data and code is permitted.

443 6. Experimental Setting/Details

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

447 Answer: [Yes]

448 Justification: We have specified all the training and test details.

449 Guidelines:

- 450 • The answer NA means that the paper does not include experiments.
- 451 • The experimental setting should be presented in the core of the paper to a level of detail
452 that is necessary to appreciate the results and make sense of them.
- 453 • The full details can be provided either with the code, in appendix, or as supplemental
454 material.

455 7. Experiment Statistical Significance

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

458 Answer: [Yes]

459 Justification: We have used Analysis of Variance to compare the means of our metrics
460 system.

461 Guidelines:

- 462 • The answer NA means that the paper does not include experiments.
- 463 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
464 dence intervals, or statistical significance tests, at least for the experiments that support
465 the main claims of the paper.

- 466 • The factors of variability that the error bars are capturing should be clearly stated (for
467 example, train/test split, initialization, random drawing of some parameter, or overall
468 run with given experimental conditions).
- 469 • The method for calculating the error bars should be explained (closed form formula,
470 call to a library function, bootstrap, etc.)
- 471 • The assumptions made should be given (e.g., Normally distributed errors).
- 472 • It should be clear whether the error bar is the standard deviation or the standard error
473 of the mean.
- 474 • It is OK to report 1-sigma error bars, but one should state it. The authors should
475 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
476 of Normality of errors is not verified.
- 477 • For asymmetric distributions, the authors should be careful not to show in tables or
478 figures symmetric error bars that would yield results that are out of range (e.g. negative
479 error rates).
- 480 • If error bars are reported in tables or plots, The authors should explain in the text how
481 they were calculated and reference the corresponding figures or tables in the text.

482 8. Experiments Compute Resources

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

486 Answer: [Yes]

487 Justification: We have added all types of computer workers and memory to better reproduce
488 our experiment.

489 Guidelines:

- 490 • The answer NA means that the paper does not include experiments.
- 491 • The paper should indicate the type of compute worker CPU or GPU, internal cluster, or
492 cloud provider, including relevant memory and storage.
- 493 • The paper should provide the amount of compute required for each of the individual
494 experimental runs and estimate the total compute.
- 495 • The paper should disclose whether the full research project required more compute
496 than the experiments reported (e.g., preliminary or failed experiments that didn't make
497 it into the paper).

498 9. Code Of Ethics

499 Question: Does the research conducted in the paper conform, in every respect, with the
500 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

501 Answer: [Yes]

502 Justification: The research conducted in the paper conform, in every respect, with the
503 NeurIPS Code of Ethics.

504 Guidelines:

- 505 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 506 • If the authors answer No, they should explain the special circumstances that require a
507 deviation from the Code of Ethics.
- 508 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
509 eration due to laws or regulations in their jurisdiction).

510 10. Broader Impacts

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

513 Answer: [Yes]

514 Justification: We have discussed how it will bring out the positive impact on the research of
515 material science with the help of LLM.

516 Guidelines:

- 517 • The answer NA means that there is no societal impact of the work performed.
- 518 • If the authors answer NA or No, they should explain why their work has no societal
- 519 impact or why the paper does not address societal impact.
- 520 • Examples of negative societal impacts include potential malicious or unintended uses
- 521 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 522 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 523 groups), privacy considerations, and security considerations.
- 524 • The conference expects that many papers will be foundational research and not tied
- 525 to particular applications, let alone deployments. However, if there is a direct path to
- 526 any negative applications, the authors should point it out. For example, it is legitimate
- 527 to point out that an improvement in the quality of generative models could be used to
- 528 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 529 that a generic algorithm for optimizing neural networks could enable people to train
- 530 models that generate Deepfakes faster.
- 531 • The authors should consider possible harms that could arise when the technology is
- 532 being used as intended and functioning correctly, harms that could arise when the
- 533 technology is being used as intended but gives incorrect results, and harms following
- 534 from (intentional or unintentional) misuse of the technology.
- 535 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 536 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 537 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 538 feedback over time, improving the efficiency and accessibility of ML).

539 11. Safeguards

540 Question: Does the paper describe safeguards that have been put in place for responsible
541 release of data or models that have a high risk for misuse (e.g., pretrained language models,
542 image generators, or scraped datasets)?

543 Answer: [NA]

544 Justification: There is no high risk of using our experiment.

545 Guidelines:

- 546 • The answer NA means that the paper poses no such risks.
- 547 • Released models that have a high risk for misuse or dual-use should be released with
- 548 necessary safeguards to allow for controlled use of the model, for example by requiring
- 549 that users adhere to usage guidelines or restrictions to access the model or implementing
- 550 safety filters.
- 551 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 552 should describe how they avoided releasing unsafe images.
- 553 • We recognize that providing effective safeguards is challenging, and many papers do
- 554 not require this, but we encourage authors to take this into account and make a best
- 555 faith effort.

556 12. Licenses for existing assets

557 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
558 the paper, properly credited and are the license and terms of use explicitly mentioned and
559 properly respected?

560 Answer: [Yes]

561 Justification: We have given proper credit.

562 Guidelines:

- 563 • The answer NA means that the paper does not use existing assets.
- 564 • The authors should cite the original paper that produced the code package or dataset.
- 565 • The authors should state which version of the asset is used and, if possible, include a
- 566 URL.
- 567 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 568 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 569 service of that source should be provided.

- 570
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

578 **13. New Assets**

579 Question: Are new assets introduced in the paper well documented and is the documentation
580 provided alongside the assets?

581 Answer: [Yes]

582 Justification: The new assets are well documented.

583 Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

592 **14. Crowdsourcing and Research with Human Subjects**

593 Question: For crowdsourcing experiments and research with human subjects, does the paper
594 include the full text of instructions given to participants and screenshots, if applicable, as
595 well as details about compensation (if any)?

596 Answer: [NA]

597 Justification: We did not do our experiments on human subjects.

598 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

607 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
608 Subjects**

609 Question: Does the paper describe potential risks incurred by study participants, whether
610 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
611 approvals (or an equivalent approval/review based on the requirements of your country or
612 institution) were obtained?

613 Answer: [NA]

614 Justification: The paper doesn't have any potential risks incurred by study participants.

615 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

621
622
623
624
625

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.