# Graphusion: A RAG Framework for Scientific Knowledge Graph Construction with a Global Perspective

### Rui Yang
Duke-NUS Medical School
Singapore
yang_rui@u.nus.edu

### Boming Yang
The University of Tokyo
Tokyo, Japan
boming@g.ecc.u-tokyo.ac.jp

### Xinjie Zhao
The University of Tokyo
Tokyo, Japan
xinjie-zhao@g.ecc.u-tokyo.ac.jp

### Fan Gao
The University of Tokyo
Tokyo, Japan
fangao0802@gmail.com

### Aosong Feng
Yale University
New Haven, CT, USA
aosong.feng@yale.edu

### Sixun Ouyang
Smartor Inc.
Shanghai, China
troy.oysx@gmail.com

### Moritz Blum
Bielefeld University
Bielefeld, Germany
mblum@techfak.uni-bielefeld.de

### Tianwei She
Smartor Inc.
Shanghai, China
tianwei.v.she@gmail.com

### Yuang Jiang
Smartor Inc.
Shanghai, China
jiangyuang1995@gmail.com

### Freddy Lecue
INRIA
Paris, France
freddy.lecue@inria.fr

### Jinghui Lu
Smartor Inc.
Shanghai, China
liuxiangtian213@gmail.com

### Irene Li
University of Tokyo
Tokyo, Japan
ireneli@ds.itc.u-tokyo.ac.jp

## Abstract

Knowledge Graphs (KGs) are crucial in the field of artificial intelligence and are widely used in downstream tasks, such as question-answering (QA). The construction of KGs typically requires significant effort from domain experts. Large Language Models (LLMs) have recently been used for Knowledge Graph Construction (KGC). However, most existing approaches focus on a local perspective, extracting knowledge triplets from individual sentences or documents, missing a fusion process to combine the knowledge in a global KG. This work introduces Graphusion, a zero-shot KGC framework from free text. It contains three steps: in Step 1, we extract a list of seed entities using topic modeling to guide the final KG includes the most relevant entities; in Step 2, we conduct candidate triplet extraction using LLMs; in Step 3, we design the novel fusion module that provides a global view of the extracted knowledge, incorporating entity merging, conflict resolution, and novel triplet discovery. Results show that Graphusion achieves scores of 2.92 and 2.37 out of 3 for entity extraction and relation recognition, respectively. Moreover, we showcase how Graphusion could be applied to the Natural Language Processing (NLP) domain and validate it in an educational scenario. Specifically, we introduce TutorQA, a new expert-verified benchmark for QA, comprising six tasks and a total of 1,200 QA pairs. Using the Graphusion-constructed KG, we achieve a significant improvement on the benchmark, for example, a 9.2% accuracy improvement on sub-graph completion.

## Introduction

Retrieval-Augmented Generation (RAG) [19] combines the advantages of retrieval methods and generative models, which improves the accuracy and relevance of generated content [42]. For instance, given a free-text corpus and a query involving two related entities, RAG can retrieve relevant information and infer their relation. Therefore, RAG can be used to enhance the performance of various knowledge-intensive tasks [9, 41, 44]. However, the need for more structured and comprehensive knowledge integration has highlighted the importance of adopting Knowledge Graphs (KGs). KGs provide structured and interconnected representations of information, offering richer context and reasoning capabilities that further enhance retrieval-augmented methods in complex applications. [5, 21, 30]
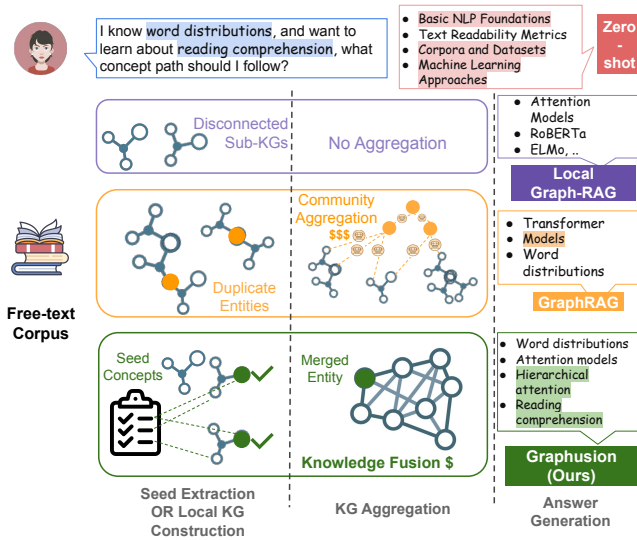
**Figure 1: Comparison of Zero-shot LLM, RAG framework, and our Graphusion framework on applying LLMs for KGC.**

Especially in the scientific domain, KGs play a crucial role when precise extraction and modeling of entities and their relations are required. The accurate selection of entities with appropriate granularity, along with precise modeling of their relations, is key to capturing the complex structure and semantics inherent in scientific knowledge [15, 20].

Existing knowledge graph construction (KGC) methods predominantly adopt a localized perspective [4, 6, 35], focusing on extracting triplets at the sentence or paragraph level. While this approach works well for shallow knowledge—such as (`people`, `belong_to`, `organization`)—it falls short in scientific domains, where a global view is essential for identifying complex, multi-layered relations between entities. Localized methods often fail to capture the comprehensive and interconnected nature of knowledge, leading to limited accuracy and completeness when triplets are derived from isolated text segments, which is particularly crucial for scientific KGs. The recent success of GraphRAG [8] highlights the value of leveraging a global KG for query-based summarization with the help of large language models (LLMs) [1]. However, despite offering enhanced contextual understanding, building and maintaining large-scale graph structures with hierarchical clustering significantly increases computational cost and complexity compared to simpler retrieval methods, limiting its application in resource-constrained environments or real-time systems requiring low-latency responses. Moreover, its effectiveness in scientific KGC, particularly when high entity granularity is required, remains unclear.

We illustrate the necessity of balancing the performance-efficiency trade-off in building the large-scale KG tailored to a user prompt in Fig. 1. A user poses a specific question from the NLP domain, requiring a learning path from the entity `word distributions` to the entity `reading comprehension`. Ideally, the model should understand the relations between these entities and effectively map out the learning path. Zero-shot LLMs might provide somewhat relevant answers but tend to be too general, as shown in the figure, offering broad entities like `Basic NLP Foundations` or
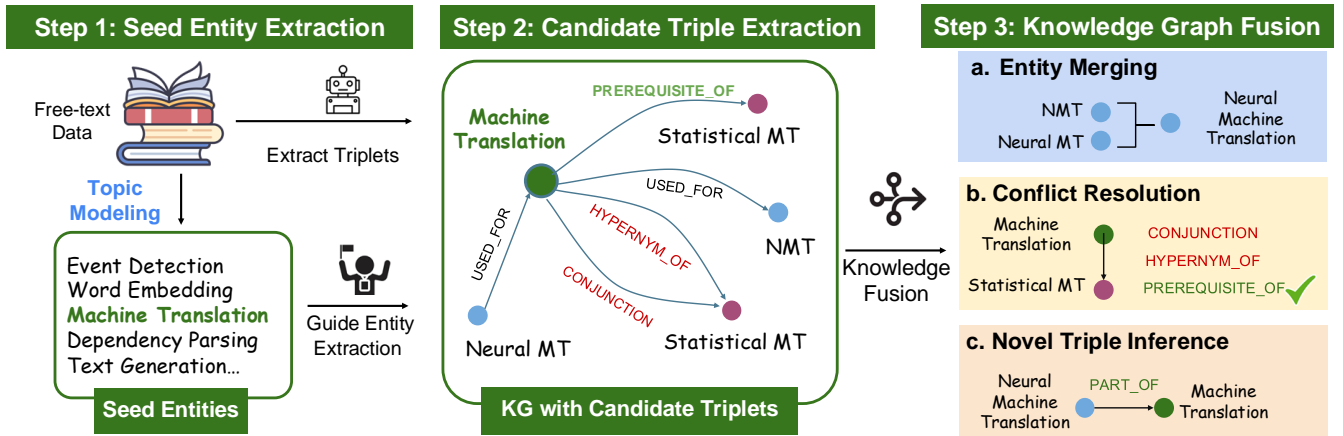
introducing confusing, inaccurately specific entities (e.g., `Corpora and Datasets`). While some RAG-based methods for building KGs focus primarily on relation extraction from limited sources via retrieval, this often results in numerous disconnected KGs or sub-graphs, we marked this method as Local Graph-RAG. We argue that a global understanding of domain knowledge is crucial and may be challenging for typical existing RAG frameworks. For instance, determining the relation between `hierarchical attention network` and `reading comprehension` may be difficult, as these entities might not appear within the same document. The model needs to extract and synthesize information from two (or more) documents to recognize that the relation is `Used_for`. GraphRAG addresses this challenge hierarchically by linking entities to gather information on a global scale. However, despite the method's high cost, it remains unclear how relation conflicts are managed. Moreover, we observe that some entities, such as `models`, may be overly broad and therefore not useful for users.

Recognizing these limitations, we propose a cost-efficient KGC approach to incorporate the global perspective and improve RAG performance on downstream question-answering (QA) tasks. Fig. 1 illustrates the differences between zero-shot LLMs and Local Graph-RAG, GraphRAG and our model. While Local Graph-RAG primarily focuses on knowledge extraction from a disconnected point of view, GraphRAG is able to summarize the information from a global view. However, our approach, Graphusion, incorporates a knowledge fusion step to integrate local knowledge into a global context directly. The core fusion step performs global merging and resolution across multiple local sub-graphs to form one large connected KG. Specifically, we leverage LLMs not only for extraction but also for critical knowledge integration, marking the first initiative to utilize LLMs for such a comprehensive merging process. Moreover, we generate a seed entity list to guide the entity extraction with better granularity. We demonstrate Graphusion's capability in knowledge graph construction and link prediction. Our results show that Graphusion achieves scores of 2.92 and 2.37 out of 3 for entity extraction and relation recognition, respectively. Furthermore, we show that a simplified link prediction prompt outperforms supervised learning baselines, achieving a 3% higher F1 score. Then, to further show the power of the constructed KG, we showcase how it could be applied to the NLP domain, and most importantly, we validate it in an educational scenario with complex QA tasks. [1]

## Related Work

**Knowledge Graph Construction** KGC aims to create a structured representation of knowledge in the form of a KG. A KG can be created from various sources, such as databases or texts. In our work, we focus on KGC from natural language resources. Research on KGs spans various domains, including medical, legal, and more [2, 16, 18, 22, 38]. Typically, KGC from text involves several methods such as entity extraction and link prediction [27, 33], with a significant focus on supervised learning. Recently, LLMs have been utilized in KGC relying on their powerful zero-shot capabilities [6, 24, 45, 47]. Although relevant works propose pipelines for extracting knowledge, they often remain limited to localized views, such as extracting triplets from the sentence or paragraph level. In

---

[1] https://github.com/IreneZihuiLi/Graphusion

**Figure 2: Graphusion framework illustration. Graphusion consists of 3 steps: S1 Seed Entity Generation, S2 Candidate Triplet Extraction and S3 Knowledge Graph Fusion.**

contrast, our work focuses on shifting from a local perspective to a global one, aiming to generate a more comprehensive KG. Approaches such as GraphRAG [8] which uses graph indexing and community detection to generate query-focused summaries, effectively answer global questions. However, GraphRAG focuses less on the detailed steps of graph construction, such as entity resolution and relation inference. In contrast, this work places greater emphasis on the process of global KGC, including entity merging, conflict resolution, and novel triplet discovery, thereby achieving a more comprehensive and consistent knowledge representation.

**Educational Question Answering** This work also falls within the scope of applications for educational question answering. Modern NLP and Artificial Intelligence (AI) techniques have been applied to a wide range of applications, with education being a significant area. For instance, various tools have been developed focusing on writing assistance, language study, automatic grading, and quiz generation [10, 26, 34, 43, 46]. Moreover, in educational scenarios, providing responses to students still requires considerable effort, as the questions often demand a high degree of relevance to the study materials and strong domain knowledge. Consequently, many studies have concentrated on developing automatic QA models [13, 48], which tackle a range of queries, from logistical to knowledge-based questions. In this work, we integrate a free-text constructed KG for various QA tasks in NLP education.

## Graphusion: Zero-shot Knowledge Graph Construction

We now introduce our Graphusion framework for constructing scientific KGs, shown in Fig. 2. Our approach addresses three key challenges in zero-shot KGC: 1) the input consists of free text rather than a predefined list of entities; 2) the relations encompass multiple types, and conflicts may exist among them; and 3) the output is not a single binary label but a list of triplets, making evaluation more challenging.

**Problem Definition** A $KG$ is defined as a set of triplets $KG = \{(h_i, r_i, t_i) \mid h_i, t_i \in E, r_i \in R, i = 1, 2, \ldots, n\}$, where $E$ is the

set of entities, $R$ is the set of possible relations, and $n$ is the total number of triplets in the $KG$. The task of zero-shot KGC involves taking a set of free text $T$ and generating a list of triplets $(h, r, t)$ spanning a KG. Optionally, there is an expert-annotated KG, $G_E$, as input, in order to provide existing knowledge. In our setting, the number of triplets of $KG$ is much larger than $G_E$. We select the domain to be NLP, so the entities are limited to NLP entities, with other entity types such as people, and organizations not being our focus. Referring to previous works [27], we define 7 relations types: `Prerequisite_of`, `Used_for`, `Compare`, `Conjunction`, `Hyponym_of`, `Evaluate_for` and `Part_of`. We will now describe the three steps of the pipeline in detail.

## Step 1: Seed Entity Generation

Extracting domain-specific entities using LLMs under a zero-shot setting is highly challenging due to the absence of predefined entity lists. This process is not only resource-intensive but also tends to generate a large number of irrelevant entities, or entities with a bad granularity, thereby compromising the quality of extraction. To address these issues, we adopt a seed entity generation approach for efficiently extracting in-domain entities from free text [17]. Specifically, we utilize BERTopic [11] for topic modeling to identify representative entities for each topic. These representative entities serve as seed entities, denoted as $Q$. The initialized seed entities ensure high relevance in entity extraction and provide certain precision for subsequent triplet extraction.

## Step 2: Candidate Triplet Extraction

These seed entities obtained from Step 1 would guide us to conduct entity extraction. In Step 2, we begin extracting candidate triplets from the free text. Each time, we input an entity $q \in Q$ ({query}) as the query entity and retrieve documents containing this entity ({context}) through information retrieval. Our goal is to extract any potential triplet that includes this query entity. To achieve this, we design a Chain-of-Thought (CoT) [39] prompt. We first instruct the LLMs to extract in-domain entities, and then identify the possible

relations between those entities and $q$. Then, we ask LLMs to discover novel triplets, even if $q$ is not initially included. This approach ensures that the seed entities play a leading role in guiding the extraction of in-domain entities. Meanwhile, the candidate triplets will encompass novel entities. We design the **Extraction Prompt** to be the following:

```
Given a context {context} and a query
entity {query}, do the following:

1. Extract the query entity and
   in-domain entities from the context,
   which should be fine-grained...
2. Determine the relations between
   the query entity and the extracted
   entities, in a triplet format:
   (<head entity>, <relation>, <tail entity
   >)...
   {Relation Definition}
3. Please note some relations are
   strictly directional...
4. You can also extract triplets from
   the extracted entities, and the
   query entity may not be necessary
   in the triplets.
```

After processing all the queries from the seed entity list, we save all the candidate triplets. We denote this zero-shot constructed KG by the LLM as $\mathcal{ZS} - \mathcal{KG}$.

## Step 3: Knowledge Graph Fusion

The triplets extracted in the previous step provide a local view rather than a global perspective of each query entity. Due to the limitations of context length, achieving a global view is challenging. Additionally, the relations extracted between two entities can be conflicting, diverse, or incorrect, such as (neural summarization methods, Used-for, abstractive summarization) and (neural summarization methods, Hyponym-of, abstractive summarization). To address the aforementioned challenge, we propose the fusion step. This approach helps reconcile conflicting relations, integrate diverse or incorrect relations effectively, and ultimately provides a global understanding of an entity pair. Specifically, for each query entity $q$, we first query from $\mathcal{ZS} - \mathcal{KG}$, and obtain a sub-graph that contains $q$:

$$\text{LLM-KG} = \{(h, r, t) \in \mathcal{ZS}\text{-}\mathcal{KG} \mid h = q \text{ or } t = q\}.$$

Optionally, if there is an expert-annotated KG $G_E$ available, we will also query a sub-graph, marked as $\mathcal{E} - \mathcal{G}$. Moreover, we conduct a dynamic retrieval of $q$ again from the free text ({background}), to help LLMs to have a better understanding on how to resolve the conflicted triplets. This key fusion step focuses on three parts: 1) entity merging: merge semantically similar entities, i.e., NMT vs neural machine translation; 2) conflict resolution: for each entity pair, resolve any conflicts and choose the best one; and 3) novel triplet inference: propose new triplet from the background text. We utilize the following **Fusion Prompt**:

```
Please fuse two sub-knowledge graphs
about the entity: {entity}.

Graph 1: {LLM-KG}
Graph 2: {E-G}

Rules for Fusing the Graphs:
1. Union the entities and edges.
2. If two entities are similar, or
   refer to the same entity, merge
   them into one entity, keeping the
   one that is meaningful or specific.
3. Only one relation is allowed between
   two entities. If a conflict exists,
   read the ### Background to help
   you keep the correct relation...
4. Once step 3 is done, consider every
   possible entity pair not covered in
   step 2. For example, take an entity
   from Graph 1, and match it with a
   entity from Graph 2. Then, please
   refer to ### Background to summarize
   new triplets.

### Background:
{background}

{Relation Definition}
```

## Experiments on Knowledge Graph Construction

In these experiments, we investigated the general capabilities of Graphusion for KGC.

**Dataset** To conduct scientific KGC, we need a large-scale free-text corpus to serve as the knowledge source. We collect the proceedings papers from the ACL conference[2] spanning 2017-2023, which includes a total of 4,605 valid papers. Considering that abstracts provide high-density, noise-free information and save computational resources, we perform topic modeling and KGC on the paper abstracts.

**Implementation** We implement Graphusion on top of four settings with different LLMs: LLaMa3-70b[3], GPT-3.5, GPT-4 and GPT-4o. Additionally, we compare with multiple baselines, including zero-shot (GPT-4o zs) and RAG (GPT-4o RAG). We query what is the relation of a node pair as the prompt, and the zero-shot setting answers directly, while the RAG one answers with retrieved results from the same data with Graphusion.

**Baseline** We compare with a local graph model (GPT-4o Local), which equals to the Graphusion model without the fusion step (Step 3). Note that it is challenging to evaluate the entity quality, as a simple prompt will generate out-of-domain nodes, so we focus on comparing the relation quality. We also compare with the GraphRAG framework. Like Graphusion, we first tune the prompt with zero-shot and CoT settings for GraphRAG' entity/relation extraction. We define the relation types within the prompt. We also utilize GraphRAG's prompt auto-tuning ability to create community reports to adapt the generated knowledge graph to the NLP domain. The full manual promo tuning can be found in the Appendix. We then

feed the collected abstracts into GraphRAG to build the indexing pipelines. Specifically, we employ GPT-4o as the base LLM. In the query phase, we ask GraphRAG the relation between the given entity pairs.

**Evaluation Metrics** The automatic evaluation of our scientific KGC approach is challenging, due to the lack of ground truth graphs matching our setting. Therefore, we conduct a human evaluation of the constructed KG. For each model, we randomly sample 100 triplets and ask experts to assess both *entity quality* and *relation quality*, providing ratings on a scale from 1 (bad) to 3 (good). Entity quality measures the relevance and specificity of the extracted entities, while relation quality evaluates the logical accuracy of the relation between entities. We provide the annotators with the following guidelines:

**1. Entity Quality** *Excellent (3 points)*: Both entities are highly relevant and specific to the domain. At an appropriate level of detail, neither too broad nor too specific. For example, an entity could be introduced by a lecture slide page, or a whole lecture, or possibly have a Wikipedia page. *Acceptable (2 points)*: Entity is somewhat relevant, or granularity is acceptable. *Poor (1 point)*: Entity is at an inappropriate level of detail, too broad or too specific.

**2. Relation Quality** *Correct (3 points)*: The relation logically and accurately describes the relation between the head and tail entities. *Somewhat Correct (2 points)*: The relation is acceptable but has minor inaccuracies or there might be another better or correct answer. *Incorrect (1 point)*: The relation does not logically describe the relation between the entities.
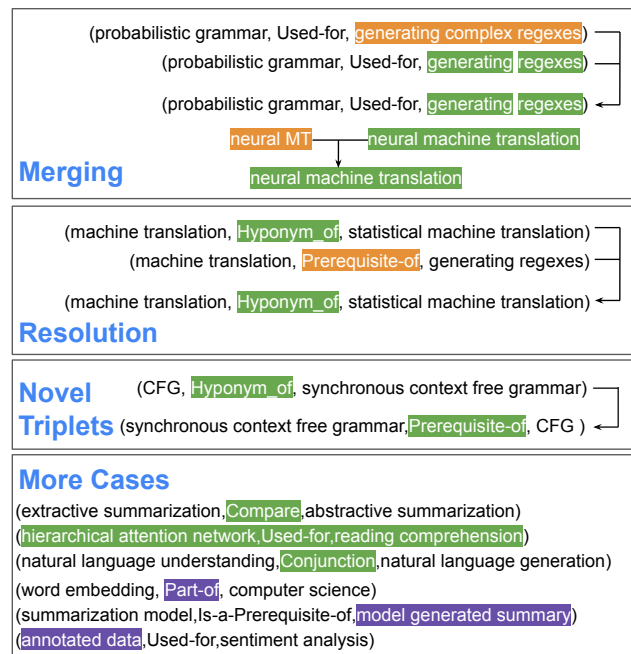
Additionally, we calculate the Inter-Annotator Agreement (IAA) between the two experts using the Kappa score.

**Results** Tab. 1 shows the ratings and the experts' consistency scores. Overall, the rating for entity surpasses relation, demonstrating the challenge of relation extraction. Among all the methods tested, Graphusion with GPT-4o achieves the highest performance in both entity and relation ratings. Notably, when the fusion step is omitted, performance drops significantly from 2.37 to 2.08, demonstrating the crucial role of the fusion step in enhancing relation quality within Graphusion. Notably, the performance of GraphRAG is suboptimal, partly due to modifications we made to better align it with our evaluation. Furthermore, our observations indicate that it often defaults to predicting a single relation type when uncertain (e.g., 40 `Part_Of` relations out of 100 cases). Additionally, the high consistency score among the experts indicates the reliability of the expert evaluation.

**Case Study: Fusion** In Fig 3, we present case studies from our Graphusion framework using GPT-4o. Our fusion step merges similar entities (`neural MT` and `neural machine translation`) and also resolves relational conflicts (`Prerequisite_of` and `Hyponym_of`). Additionally, it can infer novel triplets absent from the input. We highlight both positive and negative outputs. For instance, it correctly identifies the use of a technique for a task (`hierarchical attention network, Used_for, reading comprehension`). However, it may output less accurate triplets in entity recognition, such as entities with poor granularity (`annotated data`, `model generated summary`) and identifying very far relations (`word embedding` being categorized as part of `computer science`).

| Model | Entity | | Relation | |
|---|---|---|---|---|
| | **Rating** | **Kappa** | **Rating** | **Kappa** |
| GPT-4o zs | - | - | $2.28_{\pm 0.88}$ | 0.68 |
| GPT-4o RAG | - | - | $2.28_{\pm 0.87}$ | 0.66 |
| GPT-4o Local | - | - | $2.08_{\pm 0.86}$ | 0.59 |
| GraphRAG | - | - | $2.09_{\pm 0.70}$ | 0.56 |
| *Graphusion* | | | | |
| LLaMA | $2.83_{\pm 0.47}$ | 0.63 | $1.82_{\pm 0.81}$ | 0.51 |
| GPT-3.5 | $2.90_{\pm 0.38}$ | 0.48 | $2.14_{\pm 0.83}$ | 0.67 |
| GPT-4 | $2.84_{\pm 0.50}$ | 0.68 | $2.36_{\pm 0.81}$ | 0.65 |
| **GPT-4o** | $\mathbf{2.92_{\pm 0.32}}$ | 0.65 | $\mathbf{2.37_{\pm 0.82}}$ | 0.67 |

**Table 1: Rating for the quality of entity and relation, and IAA score for the expert evaluation.**



**Figure 3: Case studies for Graphusion on the GPT-4o model: Correct parts are highlighted in green, resolved and merged parts in orange, and less accurate parts in purple.**

## Experiments on Link Prediction

While the task of KGC is to generate a list of triplets, including entities and their corresponding relations, we also evaluate a subtask: focusing solely on Link Prediction (LP) for pre-defined entity pairs and a single relation type, $r$ =`Prerequisite_of`. Specifically, given an entity pair $(A, B)$, the task is to determine if a relation $r$ exists between two given entities. For instance, to learn the entity of `POS Tagging`, one must first understand the `Viterbi Algorithm`. Initially, a predefined set of entities $C$ is given.

We then design an **LP Prompt** to solve the task. The core part is to provide the domain name, the definition and description of the dependency relation to be predicted, and the query entities. Meanwhile, we explore whether additional information, such as entity definitions from Wikipedia and neighboring entities from training data (when

| Method | NLP | | CV | | BIO | | Overall | |
|---|---|---|---|---|---|---|---|---|
| | **Acc** | **F1** | **Acc** | **F1** | **Acc** | **F1** | **Acc** | **F1** |
| *Supervised Baselines* | | | | | | | | |
| P2V[3] | 0.6369 | 0.5961 | 0.7642 | 0.7570 | 0.7200 | 0.7367 | 0.7070 | 0.6966 |
| BERT[7] | 0.7088 | 0.6963 | 0.7572 | 0.7495 | 0.7067 | 0.7189 | 0.7242 | 0.7216 |
| DeepWalk[31] | 0.6292 | 0.5860 | 0.7988 | 0.7910 | 0.7911 | 0.8079 | 0.7397 | 0.7283 |
| Node2vec[12] | 0.6209 | 0.6181 | 0.8197 | 0.8172 | 0.7956 | 0.8060 | 0.7454 | 0.7471 |
| *Zero-shot (zs)* | | | | | | | | |
| LLaMA | 0.6058 | 0.6937 | 0.6092 | 0.6989 | 0.6261 | 0.6957 | 0.6137 | 0.6961 |
| GPT-3.5 | 0.6123 | 0.7139 | 0.6667 | 0.7271 | 0.6696 | 0.6801 | 0.6495 | 0.7070 |
| GPT-4 | 0.7639 | 0.7946 | **0.7391** | **0.7629** | **0.7348** | **0.7737** | **0.7459** | **0.7771** |
| *Zero-shot + RAG* | | | | | | | | |
| GPT-3.5 | 0.7587 | 0.7793 | 0.6828 | 0.7123 | 0.6870 | 0.7006 | 0.7095 | 0.7307 |
| GPT-4 | **0.7755** | **0.7958** | 0.7230 | 0.7441 | 0.7174 | 0.7200 | 0.7386 | 0.7533 |

**Table 2: Link prediction results across all domains on the LectureBankCD test set: We present accuracy (Acc) and F1 scores. Bolded figures indicate the best performance in the zero-shot setting, while underlined scores represent the highest achievements in the supervised setting. We apply LLaMA2-70b for all for this task.**

| Dataset | Domain | Answer Type | With KG | Collection |
|---|---|---|---|---|
| CBT [14] | Open-Domain | Multiple Choice | No | Automated |
| LectureBankCD [20] | NLP, CV, BIO | Binary | Yes | Expert-verified |
| FairytaleQA [40] | Open-Domain | Open-ended | No | Expert-verified |
| ChaTa [13] | CS | Open-ended | No | Students |
| ExpertQA [28] | Science | Open-ended | No | Expert-verified |
| SyllabusQA [29] | Multiple | Open-ended | No | Course syllabi |
| TutorQA (this work) | NLP | Open-ended, Entity List, Binary | Yes | Expert-verified |

**Table 3: Comparison with other similar benchmarks: Educational or general QA benchmarks.**

available), would be beneficial. We provide detailed prompts in the Appendix.

We conduct a comprehensive evaluation on a scientific benchmark, LectureBankCD [23], which contains entity pairs and the prerequisite labels among three domains: NLP, computer vision (CV), and bioinformatics (BIO). There are 1551, 871 and 234 entity pairs, respectively. We follow the same setting and training/testing split provided by the authors [23] and report the accuracy and F1 score, shown in Tab. 2. We compare supervised baselines, zero-shot link prediction, and zero-shot approaches using RAG models. Specifically, the RAG data predominantly consists of NLP-related content, which explains the lack of noticeable improvement in the CV and BIO domains when using RAG. Overall, the LLM method outperforms traditional supervised baselines, suggesting that LLMs have the potential to achieve higher quality in knowledge graph construction, particularly in relation prediction.

## TutorQA: A Scientific Knowledge Graph QA Benchmark

We aim to evaluate the practical usefulness of the Graphusion-constructed KG from an educational perspective. In NLP classes, students often have specific questions that require answers grounded in NLP domain knowledge, rather than general or logistical queries

related to the course. To address this need, we introduce the TutorQA benchmark, a QA dataset designed for scientific KG QA.

TutorQA consists of six categories, encompassing a total of 1,200 QA pairs that have been validated by human experts, simulating questions typically encountered in classes. These questions extend beyond simple syllabus inquiries, covering more complex and challenging topics that require KG reasoning, along with proficiency in text comprehension and question answering. We list some similar benchmarks in Tab 3. While numerous open-domain QA benchmarks exist, our focus has been primarily on those within the scientific domain and tailored for college-level education, aligning with our objective to compare with benchmarks that can emulate a learning scenario. Among those, TutorQA is distinguished by its diversity in answer types and features expert-verified questions, ensuring a high standard of quality and relevance.

**TutorQA Tasks** We design different difficulty levels of the questions and divide them into 6 tasks. We summarize the tasks and provide example data in Fig 4. More data statistics and information can be found in the supplementary materials.

*Task 1: Relation Judgment* The task is to assess whether a given triplet, which connects two entities with a relation, is accurate.

*Task 2: Prerequisite Prediction* The task helps students by mapping out the key entities they need to learn first to understand a complex target topic.

**Task 1: Relation Judgment**

**Question:** In the field of Natural Language Processing, I have come across the concepts of Penn Treebank and first-order logic. Considering the relation of Hyponym-Of, which establishes a hierarchical relationship where one entity is a more specific version or subtype of another, would it be accurate to say that the concept "Penn Treebank" is a hyponym of "first-order logic"?
**Answer:** No.
**Evaluation:** Accuracy

**Task 2: Prerequisite Prediction**

**Question:** In the domain of Natural Language Processing, I want to learn about Meta-Learning, what concepts should I learn first?
**Answer:** probabilities, optimization, machine learning resources, loss function
**Evaluation:** Similarity Score

**Task 3: Path Searching**

**Question:** In the domain of Natural Language Processing, I know about the concept of optimization, now I want to learn about the concept of neural language modeling, what concept path should I follow?
**Answer:** optimization, machine learning resources, semi-supervised learning, neural networks, neural language modeling
**Evaluation:** Similarity Score

**Task 4: Subgraph Completion**

**Question:** Given the following triplets constituting a sub-graph, please infer the relationship between "story ending generation" and "natural language generation."
Triplets: story ending generation - Is-a-Prerequisite-of - sentiment control; sentence generation - Is-a-Prerequisite-of - NLG; natural language generation - Conjunction - natural language understanding
Relationships Types: Compare, Part-of, Hyponym-Of ...
**Answer:** Hyponym-Of
**Evaluation:** Accuracy

**Task 5: Clustering**

**Question:** Given the concept PCA, can you provide some similar concepts? Please provide some similar concepts.
**Answer:** Canonical Correlation Analysis, matrix factorization, linear discriminant analysis, singular value decomposition; maximum likelihood estimation.
**Evaluation:** Hit Rate

**Task 6: Idea Hamster**

**Question:** I already know about sentiment analysis, social media analysis, sentence simplification, text summarization, citation networks. In the domain of Natural Language Processing, what potential project can I work on? Give me a possible idea. Show me the title and project description.
**Answer:** (open ended)

**Figure 4: TutorQA tasks: We present a sample data instance and the corresponding evaluation metric for each task. Note: Task 6 involves open-ended answers, which are evaluated through human assessment.**

*Task 3: Path Searching* This task helps students identify a sequence of intermediary entities needed to understand a new target entity by charting a path from the graph.

*Task 4: Sub-graph Completion* The task involves expanding the KG by identifying hidden associations between entities in a sub-graph.

*Task 5: Similar Entities* The task requires identifying entities linked to a central idea to deepen understanding and enhance learning, aiding in the creation of interconnected curriculums.

*Task 6: Idea Hamster* The task prompts participants to develop project proposals by applying learned entities to real-world contexts, providing examples and outcomes to fuel creativity.

**Scientific Knowledge Graph Question Answering** To address TutorQA tasks, we first utilize the Graphusion framework to construct an NLP KG. Then we design a framework for the interaction between the LLM and the graph, which includes two steps: command query and answer generation. In the command query stage, an LLM independently generates commands to query the graph upon receiving the query, thereby retrieving relevant paths. During the answer generation phase, these paths are provided to the LLM as contextual prompts, enabling it to perform QA.

**Evaluation Metrics** *Accuracy* We report the accuracy score for Task 1 and Task 4, as they are binary classification tasks.

*Similarity score* For Tasks 2 and 3, the references consist of a list of entities. Generally, LLMs demonstrate creativity by answering with novel entities, which are often composed of more contemporary and fresh words, even though they might not exactly match the words in the graph. Consequently, conventional evaluation metrics like keyword matching are unsuitable for these tasks. To address this, we propose the *similarity score*. This metric calculates the semantic

similarity between the entities in the predicted list $C_{pred}$ and the ground truth list $C_{gold}$. Specifically, as shown in Eq 1, for an entity $m$ from the predicted list, and an entity $n$ from the ground truth list, we calculate the cosine similarity between their embeddings achieved from pre-trained BERT model [7]. We then average these similarity scores to obtain the similarity score.

$$Score = \frac{\sum_{m \in C_{pred}} \sum_{n \in C_{gold}} sim(m, n)}{|C_{pred}| \times |C_{gold}|} \qquad (1)$$

By averaging the similarity scores, the final score provides a comprehensive measure of the overall semantic alignment between the predicted and ground truth entities.

*Hit Rate* For Task 5, we employ the classical Hit Rate metric, expressed as a percentage. This measure exemplifies the efficiency of LLM at retrieving and presenting relevant entities in its output as compared to a provided list of target entities.

*Expert Evaluation* In Task 6, where open-ended answers are generated without gold-standard responses, we resort to expert evaluation for comparative analysis between baseline results and our model. Despite available LLM-centric metrics like G-Eval [25], the specific evaluation needs of this task warrant distinct criteria, particularly examining the persuasive and scientifically sound elements of generated project proposals. Four evaluation criteria, rated on a 1-5 scale, are employed: *Entity Relevancy*: the project's alignment with the query entities. *Entity Coverage*: the extent to which the project encompasses the query entities. *Project Convincity*: the persuasiveness and practical feasibility of the project. *Scientific Factuality*: the scientific accuracy of the information within the project.

**Experimental Results** Our analysis in Tab. 5 compares the results based on Graphusion and two baselines, including GPT-4o

| Question | I already know about *neural question answering, social media analysis, topic modeling, relation extraction, and word embedding variations.* In the domain of natural language processing, what potential project can I work on? Give me a possible idea. Show me the title and project description. |
|---|---|
| **LLaMA** | **Title:** Exploring the Impact of Emojis on Textual Emotions; **Description:** Emojis have become an integral part of online communication...The project will involve the following tasks: 1. Data Collection: We will collect a dataset of text messages...and their corresponding emotional labels (positive, negative, neutral). 2. Data Preprocessing: We will preprocess the data by removing stop words... |
| **GPT-4o** | **Title:** Multi-faceted Analysis of Social Media Discussions on Climate Change; **Description:** This project aims to ... Utilizing neural question answering, we can automatically address common inquiries people have about climate topics. Social media analysis will uncover trends, sentiment, and influential users shaping the conversation. Topic modeling will identify underlying themes in these discussions ... |
| **Ours** | **Title:** Enhanced Social Media Insights through Multilingual Aspect-Based Sentiment Analysis; **Description::** This project aims to develop a comprehensive system for analyzing social media data to extract and understand user sentiments across multiple languages. By leveraging neural question answering, topic modeling, and relation extraction, the system will identify key aspects and sentiments from posts. Dependency parsing and event extraction will help in understanding syntactic and event-based relations... Word embedding variations and multilingual word embeddings will ensure accurate sentiment analysis across languages... |

**Table 4: Case study on TutorQA Task 6: LLaMA, GPT-4o, and our pipeline with constructed KG from Graphusion (part). We highlight the relevant entities, irrelevant entities, entity details and expanded relevant entities.**

| Setting | T1 | T2 | T3 | T4 | T5 |
|---|---|---|---|---|---|
| GPT4o zs | 69.20 | 64.42 | 66.61 | 44.00 | 11.45 |
| GPT4o RAG | 64.40 | 65.06 | 69.31 | 40.80 | 10.02 |
| Ours | **92.00** | **80.29** | **77.85** | **50.00** | **15.65** |

**(a) Evaluation on Tasks 1-5. T1, T4: accuracy; T2, T3: similarity score; T5: hit rate.**

| Model | Relevancy | Coverage | Convincity | Factuality |
|---|---|---|---|---|
| GPT4o zs | 4.75 | 4.84 | 4.38 | 4.63 |
| GPT4o RAG | 4.73 | 4.71 | 4.58 | 4.71 |
| Ours | **4.85** | **4.91** | **4.72** | **4.77** |

**(b) Expert evaluation on Task 6.**

**Table 5: Results for TutorQA evaluations across various tasks.**

zero-shot (zs) and GPT-4o with RAG (RAG). Our method with the Graphusion constructed KG shows significant improvements across Tasks 1 to 6 over the baselines. Specifically, Task 6 is evaluated by two NLP experts, with a Kappa score of 0.67, which suggests substantial agreement. The results indicate that our pipeline exhibits a marginally superior performance, particularly in the expert evaluation of *Convincity* and *Factuality*. This suggests that our method might be better at generating content that is not only factually accurate but also presents it in a more persuasive way to the reader. Compared to the base RAG framework, the Graphusion-generated KGs lead to better performance, particularly in Task 4 and 5, where a global understanding is essential. This improvement highlights the critical role of our core fusion step in addressing complex QA.

**Case Study: Task 6 (Expanded relevant entities in the answer)** To further understand how KGs could help in advanced educational scenarios, we present a case study on Task 6 in Tab. 4. The posed question incorporates five entities (highlighted in blue), with the task being to formulate a feasible project proposal. Although LLaMA offers a substantial project description, its content and relevance to the highlighted entities (marked in orange) are somewhat lacking. In contrast, GPT-4o not only references the queried entities but also provides detailed insights (highlighted in purple) on their potential utility within the project, such as the role of `neural question`

`answering`. Lastly, with Graphusion constructed KG, the model provides a more comprehensive solution, elaborating on the entities and introducing additional ones (highlighted in lavender) that come from the recovered graph, like `dependency parsing` and `event extraction`, while initially addressing the queried entities.

## Extension on Japanese Medical Data

In our exploration of extending Graphusion to Japanese medical data, we utilized a dataset comprising approximately 0.1 billion tokens collected from Japanese drug instructions through data crawling [32]. Example triplets generated by Graphusion include: (バラシクロビル錠500mg, 抑制される, 発疹) ((Valacyclovir Tablets 500mg, suppressed, rash)), (ゾビラックス錠400, 作用機序, ウイルスDNAの複製を阻害することによりウイルスの増殖を抑える) ((Zovirax Tablets 400, mechanism of action, inhibits the replication of viral DNA to suppress the proliferation of the virus)). We randomly selected several case studies and found that the generated triplets were reasonable, demonstrating that Graphusion exhibits good generalizability. However, conducting a comprehensive evaluation is challenging due to the significant human effort required; therefore, we leave this evaluation for future work.

## Conclusion

In this work, we proposed the Graphusion to construct scientific KGs from free text using LLMs. Through three key steps: seed entity generation, candidate triplet extraction, and KG Fusion, Graphusion builds KGs from a global perspective, addressing the limitations of traditional KGC methods. Additionally, we introduced the new benchmark dataset TutorQA, which encompasses 1,200 expert-verified QA pairs across six tasks. TutorQA is specifically designed for KG-based QA in the NLP educational scenario. We developed an automated pipeline that leveraged the Graphusion-constructed KG, significantly enhancing the performance on TutorQA compared to pure LLM baselines.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Kian Ahrabian, Xinwei Du, Richard Delwin Myloth, Arun Baalaaji Sankar Ananthan, and Jay Pujara. 2023. PubGraph: A Large-Scale Scientific Knowledge Graph. *arXiv preprint arXiv:2302.02231* (2023).

[3] Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. Unsupervised Statistical Machine Translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium.

[4] Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-Augmented Language Model Prompting for Zero-Shot Knowledge Graph Question Answering. *ArXiv* abs/2306.04136 (2023). https://api.semanticscholar.org/CorpusID:259095910

[5] Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, Florence, Italy, 4762–4779. https://doi.org/10.18653/v1/P19-1470

[6] Salvatore M. Carta, Alessandro Giuliani, Lee Cecilia piano, Alessandro Sebastian Podda, Livio Pompianu, and Sandro Gabriele Tiddia. 2023. Iterative Zero-Shot LLM Prompting for Knowledge Graph Construction. *ArXiv* abs/2307.01128 (2023). https://api.semanticscholar.org/CorpusID:259316469

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *North American Chapter of the Association for Computational Linguistics*. https://api.semanticscholar.org/CorpusID:52967399

[8] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From Local to Global: A Graph RAG Approach to Query-Focused Summarization. *ArXiv* abs/2404.16130 (2024). https://api.semanticscholar.org/CorpusID:269363075

[9] Fan Gao, Hang Jiang, Rui Yang, Qingcheng Zeng, Jinghui Lu, Moritz Blum, Tianwei She, Yuang Jiang, and Irene Li. 2024. Evaluating large language models on wikipedia-style survey generation. In *Findings of the Association for Computational Linguistics ACL 2024*. 5405–5418.

[10] Cristian D. González-Carrillo, Felipe Restrepo-Calle, Jhon Jairo Ramírez-Echeverry, and Fabio A. González. 2021. Automatic Grading Tool for Jupyter Notebooks in Artificial Intelligence Courses. *Sustainability* (2021). https://api.semanticscholar.org/CorpusID:243477284

[11] Maarten R. Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *ArXiv* abs/2203.05794 (2022). https://api.semanticscholar.org/CorpusID:247411231

[12] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 855–864.

[13] Yann Hicke, Anmol Agarwal, Qianou Ma, and Paul Denny. 2023. ChaTA: Towards an Intelligent Question-Answer Teaching Assistant using Open-Source LLMs. *ArXiv* abs/2311.02775 (2023). https://api.semanticscholar.org/CorpusID:265033489

[14] Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2015. The Goldilocks Principle: Reading Children's Books with Explicit Memory Representations. *CoRR* abs/1511.02301 (2015). https://api.semanticscholar.org/CorpusID:14915449

[15] Pengcheng Jiang, Cao Xiao, Adam Richard Cross, and Jimeng Sun. 2023. GraphCare: Enhancing Healthcare Predictions with Personalized Knowledge Graphs. In *The Twelfth International Conference on Learning Representations*.

[16] Aparna Kalla, R Shailesh, S. Preetha, Snehal Chandra, and Sudeepa Roy. 2023. Scientific Knowledge Graph Creation and Analysis. *2023 IEEE 8th International Conference for Convergence in Technology (I2CT)* (2023), 1–5. https://api.semanticscholar.org/CorpusID:258870236

[17] Yuhe Ke, Rui Yang, and Nan Liu. 2024. Comparing Open-Access Database and Traditional Intensive Care Studies Using Machine Learning: Bibliometric Analysis Study. *Journal of Medical Internet Research* 26 (2024), e48330.

[18] Anh Le-Tuan, Carlos Franzreb, Sonja Schimmler, and Manfred Hauswirth. 2022. Towards Building Live Open Scientific Knowledge Graphs. *Companion Proceedings of the Web Conference 2022* (2022). https://api.semanticscholar.org/CorpusID:248347985

[19] Patrick Lewis, Ethan Perez, Aleksandara Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *ArXiv* abs/2005.11401 (2020). https://api.semanticscholar.org/CorpusID:218869575

[20] Irene Li, Vanessa Yan, Tianxiao Li, Rihao Qu, and Dragomir R. Radev. 2021. Unsupervised Cross-Domain Prerequisite Chain Learning using Variational Graph Autoencoders. In *Annual Meeting of the Association for Computational Linguistics*. https://api.semanticscholar.org/CorpusID:234334083

[21] Irene Li and Boming Yang. 2023. NNKGC: Improving Knowledge Graph Completion with Node Neighborhoods. In *Proceedings of the Workshop on Deep Learning for Knowledge Graphs (DL4KG 2023) co-located with the 21th International Semantic Web Conference (ISWC 2023), Athens, November 6-10, 2023 (CEUR Workshop Proceedings, Vol. 3559)*, Mehwish Alam and Michael Cochez (Eds.). CEUR-WS.org. https://ceur-ws.org/Vol-3559/paper-6.pdf

[22] Irene Z Li, Alexander R. Fabbri, Swapnil Hingmire, and Dragomir R. Radev. 2020. R-VGAE: Relational-variational Graph Autoencoder for Unsupervised Prerequisite Chain Learning. *ArXiv* abs/2004.10610 (2020). https://api.semanticscholar.org/CorpusID:216056469

[23] Irene Z Li, Vanessa Yan, and Dragomir R. Radev. 2021. Efficient Variational Graph Autoencoders for Unsupervised Cross-domain Prerequisite Chains. *ArXiv* abs/2109.08722 (2021). https://api.semanticscholar.org/CorpusID:237571655

[24] Qian Li, Zhuo Chen, Cheng Ji, Shiqi Jiang, and Jianxin Li. 2024. LLM-based Multi-Level Knowledge Generation for Few-shot Knowledge Graph Completion. *Proceedings of the Thirty-ThirdInternational Joint Conference on Artificial Intelligence* (2024). https://api.semanticscholar.org/CorpusID:271494703

[25] Yang Liu, Dan Iter, Yichong Xu, Shuo Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment. In *Conference on Empirical Methods in Natural Language Processing*. https://api.semanticscholar.org/CorpusID:257804696

[26] Qingyu Lu, Baopu Qiu, Liang Ding, Liping Xie, and Dacheng Tao. 2023. Error Analysis Prompting Enables Human-Like Translation Evaluation in Large Language Models: A Case Study on ChatGPT. *ArXiv* abs/2303.13809 (2023). https://api.semanticscholar.org/CorpusID:257756967

[27] Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-Task Identification of Entities, Relations, and Coreference for Scientific Knowledge Graph Construction. *ArXiv* abs/1808.09602 (2018). https://api.semanticscholar.org/CorpusID:52118895

[28] Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. 2023. ExpertQA: Expert-Curated Questions and Attributed Answers. *ArXiv* abs/2309.07852 (2023). https://api.semanticscholar.org/CorpusID:261823130

[29] Andrew Lan Nigel Fernandez, Alexander Scarlatos. 2024. SyllabusQA: A Course Logistics Question Answering Dataset. *ArXiv* abs/2403.14666 (2024). https://api.semanticscholar.org/CorpusID:268667283

[30] Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying Large Language Models and Knowledge Graphs: A Roadmap. *IEEE Trans. Knowl. Data Eng.* 36, 7 (2024), 3580–3599. https://doi.org/10.1109/TKDE.2024.3352100

[31] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. DeepWalk: Online Learning of Social Representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New York, New York, USA) *(KDD '14)*. ACM, New York, NY, USA, 701–710. https://doi.org/10.1145/2623330.2623732

[32] Pengcheng Qiu, Chaoyi Wu, Xiaoman Zhang, Weixiong Lin, Haicheng Wang, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2024. Towards building multilingual language model for medicine. *Nature Communications* 15, 1 (2024), 8384.

[33] Justin T. Reese, Deepak R. Unni, Tiffany J. Callahan, Luca Cappelletti, Vida Ravanmehr, Seth Carbon, Tommaso Fontana, Hannah Blau, Nicolas Matentzoglu, Nomi L. Harris, Monica C. Munoz-Torres, Peter N. Robinson, marcin p. joachimiak, and Chris J. Mungall. 2020. KG-COVID-19: A Framework to Produce Customized Knowledge Graphs for COVID-19 Response. *Patterns* 2 (2020). https://api.semanticscholar.org/CorpusID:221191594

[34] Dominic Seyler, Mohamed Yahya, and Klaus Berberich. 2015. Generating Quiz Questions from Knowledge Graphs. *Proceedings of the 24th International Conference on World Wide Web* (2015). https://api.semanticscholar.org/CorpusID:7522972

[35] Jiawei Sheng, Shu Guo, Zhenyu Chen, Juwei Yue, Lihong Wang, and Tingwen Liu. 2022. Challenging the Assumption of Structure-based embeddings in Few- and Zero-shot Knowledge Graph Completion. In *International Conference on Language Resources and Evaluation*. https://api.semanticscholar.org/CorpusID:252376765

[36] Deshraj Yadav Taranjeet Singh. 2023. Embedchain: The Open Source RAG Framework. https://github.com/embedchain/embedchain.

[37] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).

[38] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (sep 2014), 78–85. https://doi.org/10.1145/2629489

[39] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022. Chain of Thought Prompting Elicits Reasoning in Large Language Models. *ArXiv* abs/2201.11903 (2022). https://api.semanticscholar.org/CorpusID:246411621

[40] Ying Xu, Dakuo Wang, Mo Yu, Daniel Ritchie, Bingsheng Yao, Tong-shuang Sherry Wu, Zheng Zhang, Toby Jia-Jun Li, Nora Bradford, Branda Sun, Tran Bao Hoang, Yisi Sang, Yufang Hou, Xiaojuan Ma, Diyi Yang, Nanyun Peng, Zhou Yu, and Mark Warschauer. 2022. Fantastic Questions and Where to Find Them: FairytaleQA – An Authentic Dataset for Narrative Comprehension. In *Annual Meeting of the Association for Computational Linguistics*. https://api.semanticscholar.org/CorpusID:247762948

[41] Rui Yang, Haoran Liu, Edison Marrese-Taylor, Qingcheng Zeng, Yuhe Ke, Wanxin Li, Lechao Cheng, Qingyu Chen, James Caverlee, Yutaka Matsuo, and Irene Li. 2024. KG-Rank: Enhancing Large Language Models for Medical QA with Knowledge Graphs and Ranking Techniques. In *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, Dina Demner-Fushman, Sophia Ananiadou, Makoto Miwa, Kirk Roberts, and Junichi Tsujii (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 155–166. https://doi.org/10.18653/v1/2024.bionlp-1.13

[42] Rui Yang, Yilin Ning, Emilia Keppo, Mingxuan Liu, Chuan Hong, Danielle S Bitterman, Jasmine Chiat Ling Ong, Daniel Shu Wei Ting, and Nan Liu. 2024. Retrieval-Augmented Generation for Generative Artificial Intelligence in Medicine. *arXiv preprint arXiv:2406.12449* (2024).

[43] Rui Yang, Ting Fang Tan, Wei Lu, Arun James Thirunavukarasu, Daniel Shu Wei Ting, and Nan Liu. 2023. Large language models in health care: Development, applications, and challenges. *Health Care Science* 2, 4 (2023), 255–263.

[44] Rui Yang, Qingcheng Zeng, Keen You, Yujie Qiao, Lucas Huang, Chia-Chun Hsieh, Benjamin Rosand, Jeremy Goldwasser, Amisha Dave, Tiarnan Keenan, et al. 2024. Ascle—A Python Natural Language Processing Toolkit for Medical Text Generation: Development and Evaluation Study. *Journal of Medical Internet Research* 26 (2024), e60601.

[45] Yichi Zhang, Zhuo Chen, Wen Zhang, and Hua zeng Chen. 2023. Making Large Language Models Perform Better in Knowledge Graph Completion. *ArXiv* abs/2310.06671 (2023). https://api.semanticscholar.org/CorpusID:263830580

[46] Zheng Zhang, Jie Gao, Ranjodh Singh Dhaliwal, and Toby Jia-Jun Li. 2023. VISAR: A Human-AI Argumentative Writing Assistant with Visual Programming and Rapid Draft Prototyping. *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (2023). https://api.semanticscholar.org/CorpusID:258179241

[47] Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao, Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. LLMs for Knowledge Graph Construction and Reasoning: Recent Capabilities and Future Opportunities. *ArXiv* abs/2305.13168 (2023). https://api.semanticscholar.org/CorpusID:258833039

[48] Brian Zylich, Adam Viola, Brokk Toggerson, Lara Al-Hariri, and Andrew S. Lan. 2020. Exploring Automated Question Answering Methods for Teaching Assistance. *Artificial Intelligence in Education* 12163 (2020), 610 – 622. https://api.semanticscholar.org/CorpusID:220364751

# Zero-shot Link Prediction Prompts

*LP Prompt.*

```
We have two {domain} related entities: A: {entity_1} and B: {entity_2}.

Do you think learning {entity_1} will help in understanding {entity_2}?

Hints:
1. Answer YES or NO only.
2. This is a directional relation, which means if the answer is "YES", (B, A) is
   false, but (A, B) is true.
3. Your answer will be used to create a knowledge graph.

{Additional Information}
```

*LP Prompt With Chain-of-Thought.*

```
We have two {domain} related entities: A: {entity_1} and B: {entity_2}.

Assess if learning {entity_1} is a prerequisite for understanding {entity_2}.

Employ the Chain of Thought to detail your reasoning before giving a final answer.

# Identify the Domain and entities: Clearly define A and B within their domain.
  Understand the specific content and scope of each entity.

# Analyze the Directional Relationship: Determine if knowledge of entity A is
  essential before one can fully grasp entity B. This involves considering if A
  provides foundational knowledge or skills required for understanding B.

# Evaluate Dependency: Assess whether B is dependent on A in such a way that
  without understanding A, one cannot understand B.

# Draw a Conclusion: Based on your analysis, decide if understanding A is a
  necessary prerequisite for understanding B.

# Provide a Clear Answer: After detailed reasoning, conclude with a distinct answer
  : <result>YES</result> if understanding A is a prerequisite for understanding B,
  or <result>NO</result> if it is not.
```

*Extraction Prompt.*

```
### Instruction:
You are a domain expert in natural language processing, and now you are building a
knowledge graph in this domain.

Given a context (### Content), and a query entity (### entity), do the following:

1. Extract the query entity and in-domain entities from the context, which should
   be fine-grained: could be introduced by a lecture slide page, or a whole
   lecture, or possibly to have a Wikipedia page.

2. Determine the relations between the query entity and the extracted entities, in
   a triplet format: (<head entity>, <relation>, <tail entity>). The relation
   should be functional, aiding learners in understanding the knowledge. The query
   entity can be the head entity or tail entity.

   We define 7 types of the relations:

   a) Compare: Represents a relation between two or more entities where a
      comparison is being made. For example, "A is larger than B" or "X is more
      efficient than Y."

   b) Part-of: Denotes a relation where one entity is a constituent or component of
      another. For instance, "Wheel is a part of a Car."

   c) Conjunction: Indicates a logical or semantic relation where two or more
      entities are connected to form a group or composite idea. For example, "Salt
      and Pepper."

   d) Evaluate-for: Represents an evaluative relation where one entity is assessed
      in the context of another. For example, "A tool is evaluated for its
      effectiveness."

   e) Is-a-Prerequisite-of: This dual-purpose relation implies that one entity is
      either a characteristic of another or a required precursor for another. For
      instance, "The ability to code is a prerequisite of software development."

   f) Used-for: Denotes a functional relation where one entity is utilized in
      accomplishing or facilitating the other. For example, "A hammer is used for
      driving nails."

   g) Hyponym-Of: Establishes a hierarchical relation where one entity is a more
      specific version or subtype of another. For instance, "A Sedan is a hyponym
      of a Car."

3. Please note some relations are strictly directional. For example, "A tool is
   evaluated for B" indicates (A, Evaluate-for, B), NOT (B, Evaluate-for, A).
   Among the seven relation types, only "a) Compare" and "c) Conjunction" are not
   direction-sensitive.

4. You can also extract triplets from the extracted entities, and the query entity
   may not be necessary in the triplets.

5. Your answer should ONLY contain a list of triplets, each triplet is in this
   format: (entity, relation, entity). For example: "(entity, relation, entity)
   (entity, relation, entity)." No numbering and other explanations are needed.

6. If ### Content is empty, output None.
```

*Fusion Prompt.*

```
### Instruction: You are a knowledge graph builder.
    Now please fuse two sub-knowledge graphs about the entity "{entity}".

Graph 1: {LLM-KG}    Graph 2: {E-G}

Rules for Fusing the Graphs:
1. Union the entities and edges.

2. If two entities are similar, or refer to the same entity, merge them into one
   entity, keeping he one that is meaningful or specific. For example, "lstm"
   versus "long short-term memory",  please keep "long short-term memory".

3. Only one relation is allowed between two entities. If there is a conflict, read
   the "### Background" to help you keep the correct relation. knowledge to keep the
   correct one. For example, (ROUGE, Evaluate-for, question answering model) and
   (ROUGE,Used-for , question answering model) are considered to be conflicts.

4. Once step 3 is done, consider every possible entity pair not covered in step 2.
   For example, take an entity from Graph 1, and match it from Graph 2. Then,
   please refer to "### Background" to summarize new triplets.

Hint: the relation types and their definition. You can use it to do Step 3.
We define 7 types of the relations:

    a) Compare: Represents a relation between two or more entities where a
       comparison is being made. For example, "A is larger than B" or "X is more
       efficient than Y."

    b) Part-of: Denotes a relation where one entity is a constituent or component of
       another. For instance, "Wheel is a part of a Car."

    c) Conjunction: Indicates a logical or semantic relation where two or more
       entities are connected to form a group or composite idea. For example, "Salt
       and Pepper."

    d) Evaluate-for: Represents an evaluative relation where one entity is assessed
       in the context of another. For example, "A tool is evaluated for its
       effectiveness."

    e) Is-a-Prerequisite-of: This dual-purpose relation implies that one entity is
       either a characteristic of another or a required precursor for another. For
       instance, "The ability to code is a prerequisite of software development."

    f) Used-for: Denotes a functional relation where one entity is utilized in
       accomplishing or facilitating the other. For example, "A hammer is used for
       driving nails."

    g) Hyponym-Of: Establishes a hierarchical relation where one entity is a more
       specific version or subtype of another. For instance, "A Sedan is a hyponym
       of a Car."

### Background:
{background}

### Output Instruction:
    Output the new merged data by listing the triplets. Your answer should ONLY contain
    triplets in this format: (entity, relation, entity). No other explanations or numbering
    are needed. Only triplets, no intermediate results.
```

*Link Prediction with **Doc.***

```
We have two {domain} related entities: A: {entity_1} and B: {entity_2}.

Do you think learning {entity_1} will help in understanding {entity_2}?

Hints:
1. Answer YES or NO only.
2. This is a directional relation, which means if the answer is "YES", (B, A) is
   false, but (A, B) is true.
3. Your answer will be used to create a knowledge graph.

And here are related contents to help:
{related documents concatenation}
```

*Link Prediction with **Con.***

```
We have two {domain} related entities: A: {entity_1} and B: {entity_2}.

Do you think learning {entity_1} will help in understanding {entity_2}?

Hints:
1. Answer YES or NO only.
2. This is a directional relation, which means if the answer is "YES", (B, A) is
   false, but (A, B) is true.
3. Your answer will be used to create a knowledge graph.

And here are related contents to help:

We know that {entity_1} is a prerequisite of the following entities:
{1-hop successors of entity_1 from training data};

The following entities are the prerequisites of {entity_1}:
{1-hop predecessors of entity_1 from training data}.

We know that {entity_2} is a prerequisite of the following entities:
{1-hop successors of entity_2 from training data};

The following entities are the prerequisites of {entity_2}:
{1-hop predecessors of entity_2 from training data}.
```

*Link Prediction with **Wiki.***

```
We have two {domain} related entities: A: {entity_1} and B: {entity_2}.

Do you think learning {entity_1} will help in understanding {entity_2}?

Hints:
1. Answer YES or NO only.
2. This is a directional relation, which means if the answer is "YES", (B, A) is
   false, but (A, B) is true.
3. Your answer will be used to create a knowledge graph.

And here are related contents to help:
{Wikipedia introductory paragraph of {entity_1}}
{Wikipedia introductory paragraph of {entity_2}}
```

*GraphRAG's Prompt Tuning for Entity/Relationship Extraction.*

```
-Goal-
Given a text document that is potentially relevant to this activity, first identify all the
    entities needed from the text in order to capture the information and ideas in the text.
    Next, introduce each relation concept by defining the relation, and then report all
    relationships among the identified entities according to the predefined relations. These
    predefined relations and seed entities include:

-Relation Concepts and Definitions-:
a) Compare: Represents a relation between two or more entities where a comparison is being
    made. For example, "A is larger than B" or "X is more efficient than Y."
b) Part-of: Denotes a relation where one entity is a constituent or component of another. For
     instance, "Wheel is a part of a Car."
c) Conjunction: Indicates a logical or semantic relation where two or more entities are
    connected to form a group or composite idea. For example, "Salt and Pepper."
d) Evaluate-for: Represents an evaluative relation where one entity is assessed in the
    context of another. For example, "A tool is evaluated for its effectiveness."
e) Is-a-Prerequisite-of: This dual-purpose relation implies that one entity is either a
    characteristic of another or a required precursor for another. For instance, "The ability
     to code is a prerequisite of software development."
f) Used-for: Denotes a functional relation where one entity is utilized in accomplishing or
    facilitating the other. For example, "A hammer is used for driving nails."
g) Hyponym-of: Establishes a hierarchical relation where one entity is a more specific
    version or subtype of another. For instance, "A Sedan is a hyponym of a Car."

-Steps-
1. Identify all entities: For each identified entity, extract the following information:
- entity_name: Name of the entity,
Format each entity as ("entity"{tuple_delimiter}<entity_name>)

2. Identify all relations: From the entities identified in step 1, determine the relation
    between each pair of entities based on the predefined relation concepts (Compare, Part-of
    , Conjunction, Evaluate-for, Is-a-Prerequisite-of, Used-for, and Hyponym-of). For each
    pair of related entities:
- source_entity: Name of the source entity as identified in step 1
- target_entity: Name of the target entity as identified in step 1
- relationship_type: Select the appropriate relation from the predefined relations
- relationship_strength: a numeric score indicating strength of the relationship between the
    source entity and target entity

Format each relationship as ("relationship"{tuple_delimiter}<source_entity>{tuple_delimiter}<
    target_entity>{tuple_delimiter}<relationship_type>{tuple_delimiter}<relationship_strength
    >)
Return output: Provide the list of all entities and relationships identified in steps 1 and
    2. Use {record_delimiter} as the list delimiter. When finished, output {
    completion_delimiter}.


######################
-Real Data-:
######################
text: {input_text}
######################
output:
```

## Graphusion Case Study: Entity Extraction

To demonstrate the importance of the seed entity list from Step 1, we examine a selection of random entities extracted by both GraphRAG and Graphusion, as shown in Tab. 6. All results are based on the GPT-4o backbone. GraphRAG sometimes extracts overly general terms, such as

| GraphRAG | Graphusion |
|---|---|
| mixture-of-experts technique | code-switching tasks |
| mbart | NAS-BERT |
| benchmark | linear indexed grammars |
| multilingual alignment | analyzing data |
| diffusionbert | semantic parsing |
| proposed method | bias-variance |
| methodology | few-shot learning |

**Table 6: Entity comparison: Random samples from GraphRAG and Graphusion.**

`benchmark` and `methodology`, which occur frequently in the corpus. In subsequent experiments, we will further illustrate how GraphRAG tends to extract entities with unsuitable granularity.

## Experimental Setup

In our experimental setup, we employed Hugging Face's `LLaMA-2-70b-chat-hf`[4] and `LLaMA-3-70b-chat-hf`[5] for LLaMA2 and LLaMA3 on a cluster equipped with 8 NVIDIA A100 GPUs. For GPT-3.5 and GPT-4, we used OpenAI's `gpt-3.5-turbo`, `gpt-4-1106-preview`, and `gpt-4o` APIs, respectively, each configured with a temperature setting of zero. The RAG models are implemented using Embedchain [36]. To solve TutorQA tasks, we implemented our pipeline using LangChain[6]. The total budget spent on this project, including the cost of the GPT API service, is approximately 500 USD.

## Additional Corpora Description

**TutorialBank** We obtained the most recent version of TutorialBank from the authors, which consists of 15,583 manually curated resources. This collection includes papers, blog posts, textbook chapters, and other online resources. Each resource is accompanied by metadata and a publicly accessible URL. We downloaded the resources from these URLs and performed free text extraction. Given the varied data formats such as PDF, PPTX, and HTML, we encountered some challenges during text extraction. To ensure text quality, we filtered out sentences shorter than 25 words. Ultimately, this process yielded 559,217 sentences suitable for RAG and finetuning experiments.

**NLP-Papers** We downloaded conference papers from EMNLP, ACL, and NAACL spanning the years 2021 to 2023. Following this, we utilized Grobid (https://github.com/kermitt2/grobid) for text extraction, resulting in a collection of 4,787 documents with clean text.

## Ablation Study

**Prompting Strategies** In Tab. 7, we explore the impact of different prompting strategies for entity graph recovery, comparing CoT and zero-shot prompts across both NLP and CV domains. The results indicate the introduction of CoT is not improving. We further find that CoT Prompting more frequently results in negative predictions. This finding serves as a drawback for our study, as it somewhat suppresses the performance of our system. This observation highlights the need to balance the impact of CoT on the rigor and complexity of predictions, especially in the context of graph recovery.

| Model | NLP | | CV | |
|---|---|---|---|---|
| | Acc | F1 | Acc | F1 |
| GPT-4 zs | 0.7639 | 0.7946 | 0.7391 | 0.7629 |
| GPT-4 CoT | 0.7342 | 0.6537 | 0.6122 | 0.4159 |

**Table 7: Comparison of zero-shot and CoT prompts with GPT-4: Results on NLP and CV.**

**Finetuning** We further explore the impact of finetuning on additional datasets, with results detailed in Table 8. Specifically, we utilize LLaMA2-70b [37], finetuning it on two previously mentioned datasets: TutorialBank and NLP-Papers. Both the zero-shot LLaMA and the finetuned models are employed to generate answers. As these answers are binary (`YES` or `NO`), we can calculate both the accuracy and F1 score for evaluation. However, the results indicate that finetuning does not yield positive outcomes. This can be attributed to two potential factors: 1) the poor quality of data, and 2) limited effectiveness in aiding the graph recovery task. We leave this part as the future work.
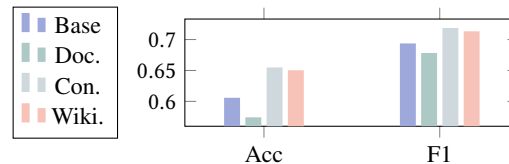
---

[4]https://huggingface.co/meta-LLaMA
[5]https://huggingface.co/meta-LLaMA/Meta-LLaMA-3-70B
[6]https://www.langchain.com/

| Dataset | Acc | F1 |
|---|---|---|
| LLaMA2-70b | **0.6058** | **0.6937** |
| TutorialBank | 0.4739 | 0.0764 |
| NLP Papers | 0.5435 | 0.6363 |

**Table 8: Comparison of the effect of finetuning: Results on NLP domain.**

## Ablation Study: RAG Data for Link Prediction

We explore the potential of external data in enhancing entity graph recovery. This is achieved by expanding the {Additional Information} part in the **LP Prompt**. We utilize LLaMa as the **Base** model, focusing on the NLP domain. We introduce three distinct settings: **Doc.**: In-domain lecture slides data as free-text; **Con.**: Adding one-hop neighboring entities from the training set as additional information related to the query entities. **Wiki.**: Incorporating the introductory paragraph of the Wikipedia page of each query entity. As illustrated in Fig 5, our findings indicate that incorporating LectureBankCD documents (Doc.) significantly diminishes performance. This decline can be attributed to the introduction of noise and excessively lengthy content, which proves challenging for the LLM to process effectively. Conversely, the inclusion of neighboring entities (Con.) markedly enhances the base model's performance. However, it relies on training data, rendering it incompatible with our primary focus on the zero-shot setting. Incorporating Wikipedia content also yields improvements and outperforms the use of LectureBankCD, likely due to higher text quality.



**Figure 5: Link Prediction Ablation Study: Comparison of models with external data.**

## TutorQA
### Benchmark Details

We show the data analysis in Tab. 9.

| Task | Question Token | | | entity Count | | | Number |
|---|---|---|---|---|---|---|---|
| | **Max** | **Min** | **Mean** | **Max** | **Min** | **Mean** | |
| T1 | 77 | 61 | 68.00 | - | - | - | 250 |
| T2 | 27 | 22 | 23.48 | 7 | 1 | 1.79 | 250 |
| T3 | 40 | 34 | 36.66 | 8 | 2 | 3.36 | 250 |
| T4 | 88 | 76 | 83.00 | - | - | - | 250 |
| T5 | 21 | 18 | 19.26 | 8 | 1 | 4.76 | 100 |
| T6 | 54 | 42 | 48.62 | - | - | - | 100 |

**Table 9: TutorQA data statistics comparison: The answers in T1 are only "True" or "False", and the answers in T4 are relations, while the answers in T6 are free text with open-ended answers.**

## GraphRAG Results

We extend the results in Tab. 5 by adding GraphRAG as a baseline, the full version of the evaluation is shown in Tab. 10. Based on the established indexing pipelines in knowledge graph construction, we utilize GraphRAG's query engine with the local search method to directly ask the questions in TutorQA. Notably, the performance of GraphRAG appears less satisfactory, which may be due to an evaluation approach that is not well-suited for GraphRAG's results. For example, in Task 5, GraphRAG produces concepts with very broad or specific terms with a bad granularity, such as *predict sentiment, emotion cause pair extraction, emotional support conversation*. This observation holds across other tasks, where achieving higher scores requires a more granular concept list. This indicates the critical importance of Step 1, which involves generating a well-defined seed concept, in the Graphusion pipeline.

| Setting | T1 | T2 | T3 | T4 | T5 |
|---|---|---|---|---|---|
| GPT4o zs | 69.20 | 64.42 | 66.61 | 44.00 | 11.45 |
| GPT4o RAG | 64.40 | 65.06 | 69.31 | 40.80 | 10.02 |
| GraphRAG | 60.40 | 64.19 | 67.45 | 42.00 | 8.96 |
| Ours | **92.00** | **80.29** | **77.85** | **50.00** | **15.65** |

(a) Evaluation on Tasks 1-5. T1, T4: accuracy; T2, T3: similarity score; T5: hit rate.

| Model | Relevancy | Coverage | Convincity | Factuality |
|---|---|---|---|---|
| GPT4o zs | 4.75 | 4.84 | 4.38 | 4.63 |
| GPT4o RAG | 4.73 | 4.71 | 4.58 | 4.71 |
| GraphRAG | 3.94 | 4.08 | 4.13 | 4.45 |
| Ours | **4.85** | **4.91** | **4.72** | **4.77** |

(b) Expert evaluation on Task 6.

**Table 10: Results for TutorQA evaluations across various tasks.**



**Figure 6: Entity counts in Task 2 and Task 3.**

## Task 2 and Task 3: case study

**Entity counts** As depicted in Fig 6, we evaluate the average number of entities created by GPT-4o zs and our Graphusion framework in the responses for Task 2 and Task 3, in which both tasks require the model to give a list of reasonable entities. The results show that without the enhancement of KG retrieved information, GPT-4o tends to mention more entities in the generated response (Task 2: 11.04, Task 3: 11.54), which might be irrelevant or broad.

| *Question* [Task2] | *In the field of Natural Language Processing, I want to learn about **multilingual model**. What entities should I learn first?* |
|---|---|
| **GPT-4o** | **Tokenization, Embeddings, Transfer Learning, Cross-lingual Transfer, Zero-shot Learning, Multilingual Corpora, Language Modeling, Fine-tuning, Evaluation Metrics, Pretrained Models** |
| **Ours** | **language models, machine translation, cross-lingual embeddings, transfer learning, tokenization, fine-tuning** |
| *Question* [Task3] | *In the field of Natural Language Processing, I know about **natural language processing intro**, now I want to learn about **t-sne**. What entity path should I follow?* |
| **GPT-4o** | **natural language processing, dimensionality reduction, t-SNE, perplexity, high-dimensional data, data visualization, machine learning** |
| **Ours** | **natural language processing intro, vector representations, t-sne** |

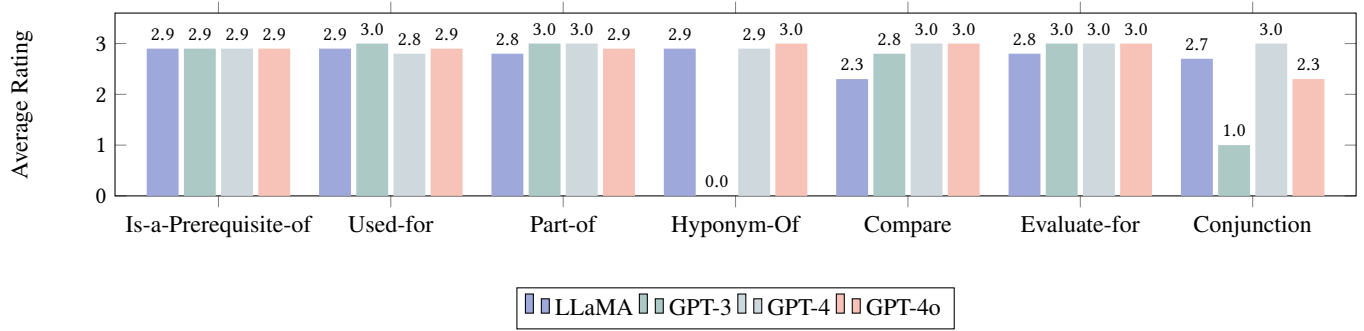**Table 11: Case study on TutorQA Task 2 and Task 3: GPT-4o, and GPT-4o-Graphusion.**

## Knowledge Graph Construction Analysis

**Average Rating** We compare expert ratings on the Graphusion KGC results produced by four models: LLaMA, GPT-3.5, GPT-4, and GPT-4o. Fig. 7 and 8 display the average ratings for entity quality and relation quality, respectively, grouped by relation type. Most types achieve an average rating of around 3 (full score) in entity quality, indicating that the extracted triplets contain good in-domain entities. In contrast, the ratings for relation quality are slightly lower. GPT-4 and GPT-4o perform better in relation prediction.
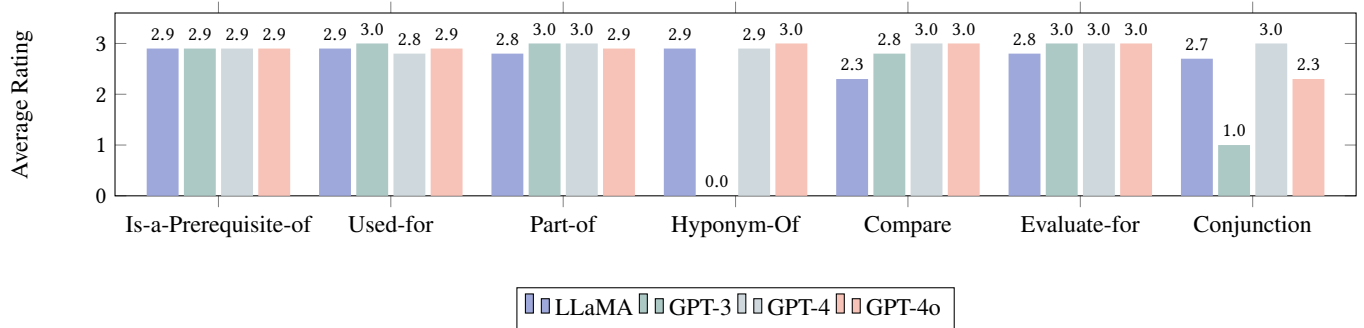
| *Question* | *Given the following edges constituting an entity subgraph, please identify and select the possible type of relationship between* **natural language generation** *and* **natural language understanding**. |
|---|---|
| **GPT-4o** | **Is-a-Prerequisite-of** |
| **Ours** | **Conjunction** |

<div align="center">

**Table 12: Case study on TutorQA Task 4: GPT-4o, and GPT-4o-Graphusion.**

</div>

<div align="center">

**Figure 7: Entity quality rating by human evaluation, grouped by relation type.**

</div>

<div align="center">

**Figure 8: Relation quality rating by human evaluation, grouped by relation type.**

</div>

**Relation Type Distribution** We then compare the Graphusion results for each relation type across the four selected base LLMs, as shown in Fig. 9. All models tend to predict `Prerequisite_of` and `Used_For` relations. The results from LLaMA show relatively even distributions across relation types, whereas the results from the GPT family do not.

**Word cloud Visualization** Finally, in Fig. 10, we present a word cloud visualization of the entities extracted by Graphusion, comparing the four base LLMs. High-frequency entities include `word embedding`, `model`, `neural network`, `language model`, and others.
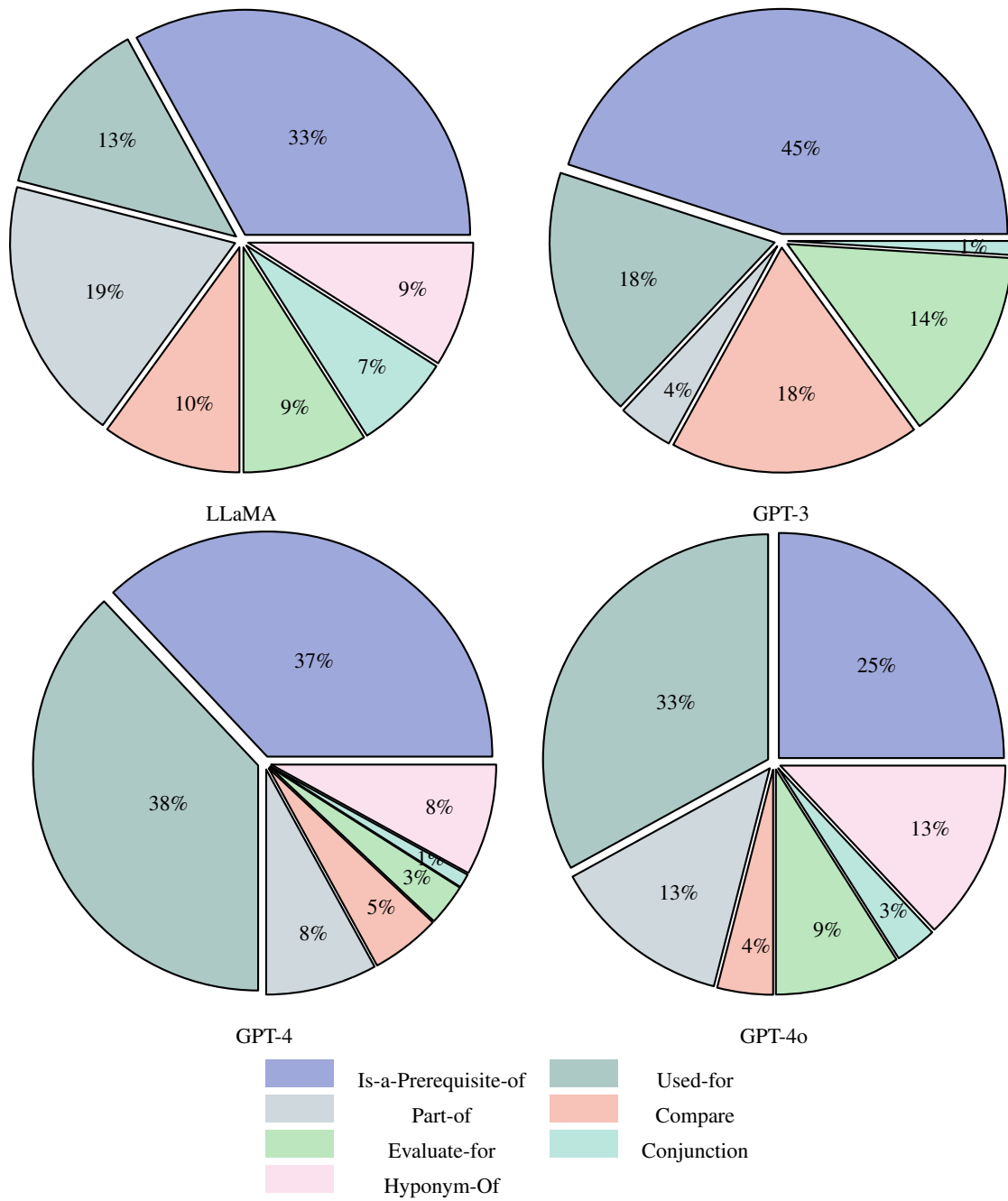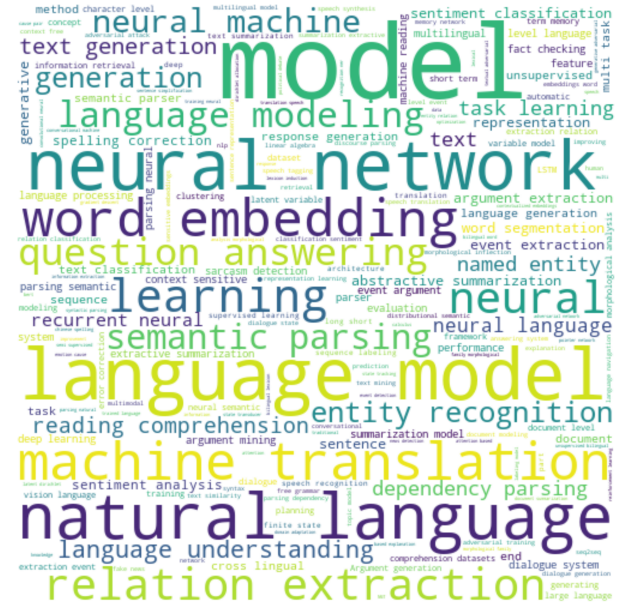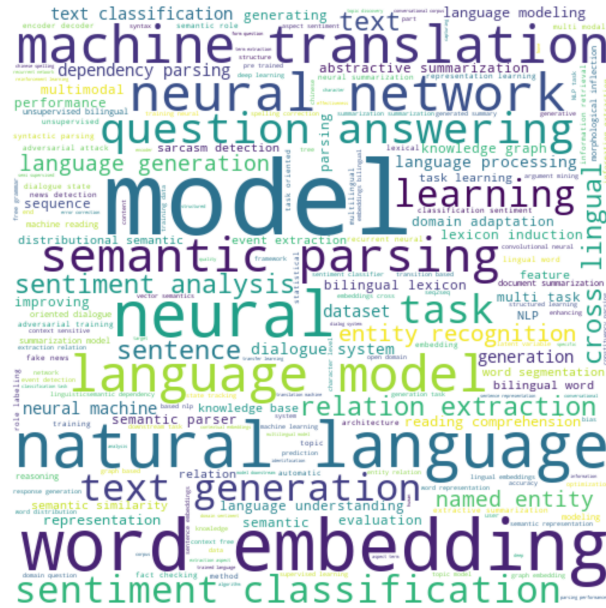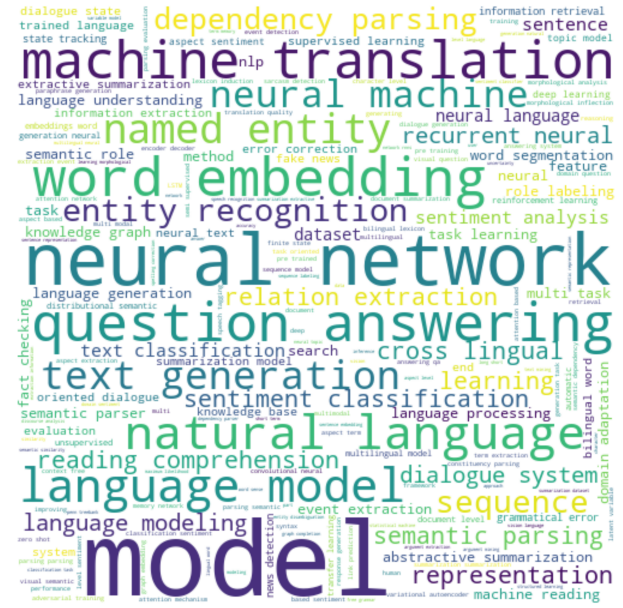
Figure 9: Relation type distribution.

**Figure 10: Word cloud visualization for extracted entities.**