# GKG-LLM: A Unified Framework for Generalized Knowledge Graph Construction

**Jian Zhang**[1,3], **Bifan Wei**[2,3], **Shihao Qi**[1,3], **Haiping Zhu**[1,3], **Jun Liu**[1,3], **Qika Lin**[4*]

[1]School of Computer Science and Technology, Xi'an Jiaotong University, Xi'an, China
[2]School of Continuing Education, Xi'an Jiaotong University, Xi'an, China
[3]Shaanxi Province Key Laboratory of Big Data Knowledge Engineering,
Xi'an Jiaotong University, Xi'an, China
[4]National University of Singapore
{zhangjian062422@stu., weibifan@, chihoq@stu., zhuhaiping@, liukeen@}xjtu.edu.cn,
linqika@nus.edu.sg

## Abstract

The construction of Generalized Knowledge Graph (GKG), including knowledge graph, event knowledge graph and commonsense knowledge graph, is fundamental for various natural language processing tasks. Current studies typically construct these types of graph separately, overlooking holistic insights and potential unification that could be beneficial in computing resources and usage perspectives. However, a key challenge in developing a unified framework for GKG is obstacles arising from task-specific differences. In this study, we propose a unified framework for constructing generalized knowledge graphs to address this challenge. First, we collect data from 15 sub-tasks in 29 datasets across the three types of graphs, categorizing them into in-sample, counter-task, and out-of-distribution (OOD) data. Then, we propose a three-stage curriculum learning fine-tuning framework, by iteratively injecting knowledge from the three types of graphs into the Large Language Models. Extensive experiments show that our proposed model improves the construction of all three graph types across in-domain, OOD and counter-task data.

## 1 Introduction

Generalized Knowledge Graph (GKG) [Krause *et al.*, 2022] includes Knowledge Graph (KG), Event Knowledge Graph (EKG) and Commonsense Knowledge Graph (CKG). The construction of GKG encompasses multiple essential tasks [Peng *et al.*, 2023], which are crucial for various applications in this field, including intelligence analysis [Pimenov *et al.*, 2023] and decision support [Lai *et al.*, 2023]. As shown in Figure 1, KGs [Lin *et al.*, 2023; Lin *et al.*, 2025a] are developed to more effectively describe concepts and relations in the physical world. The fundamental structure is <*entity, relation, entity* >, such as <Lincoln, BornIn, 1809>.
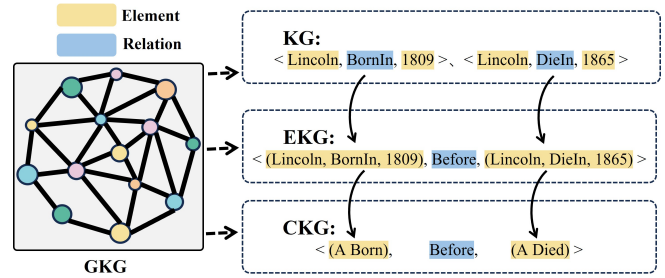
---

[*]Corresponding author



Figure 1: An illustration of several triples and graphs. The left half shows a generalized knowledge graph. The right half includes specific examples of triples from KG, EKG, CKG and demonstrates their progressive relationship.

With ongoing research, EKGs are introduced to study the dynamic progression of events. It is organized in the triplet format <*event, relation, event* >, as illustrated by <(Lincoln, BornIn, 1809), Before, (Lincoln, diedIn, 1865) >. The further generalization of event graphs has led to the development of CKG, which abstractly represent general relational patterns in the form of <*commonsense, relation, commonsense*>. For instance, <(A born), Before, (A died)>is also organized in a triplet format. In summary, KG, EKG, and CKG are all organized in the basic form of <**element, relation, element** >.

Overall, constructing the three types of graphs separately requires substantial resources, while using a unified framework for their construction improves parameter efficiency. Additionally, from a usage perspective, the knowledge contained in KGs facilitates the construction of both EKGs and CKGs. For example, a method leveraging hierarchical KGs to enhance the accuracy and effectiveness of biomedical event extraction is proposed by [Huang *et al.*, 2020]. Similarly, for knowledge graphs aiding text classification in the construction of CKGs, KG-MTT-BERT [He *et al.*, 2022] is introduced to enhance BERT with KGs for multi-type medical text classification.

Naturally, we abstract a new task to build a unified framework for constructing GKG, in order to empower these foun-

dational triples extraction tasks. However, a key challenge in this task is the obstacles arising from task-specific differences. The construction of different types of graph involves a wide variety of diverse sub-tasks. Specifically, as illustrated in Figure 2, the construction of KG includes sub-tasks such as sentence-level relation extraction [Wadhwa *et al.*, 2023], document-level relation extraction [Ma *et al.*, 2023] and joint entity and relation extraction [Sui *et al.*, 2023]. The construction of EKG involves sub-tasks such as sentence-level event detection [Hettiarachchi *et al.*, 2023], document-level argument extraction [Zhang *et al.*, 2024], and event temporal relation extraction [Chan *et al.*, 2024]. While the construction of CKG includes sub-tasks such as abstract generation [Gao *et al.*, 2023] and language inference [Gubelmann *et al.*, 2024]. The abbreviations and introduction of the task can be found in Appendix F. These tasks differ in several ways, with the primary distinctions lying in their definitions and content. For instance, sentence-level relation extraction involves extracting the relationship between two entities from a single sentence, whereas abstract generation involves extracting an abstract from an entire article. Differences between these tasks have created obstacles to building a unified framework for constructing GKG.

Thanks to the emergence of Large Language Models(LLMs), such as GPT4 [Achiam *et al.*, 2023] and LlaMA-3 [Dubey *et al.*, 2024], the realization of this new unified task has become possible. The standardized input-output format of LLMs unifies these sub-tasks from a structural perspective. To this end, we propose a three-stage curriculum learning tuning framework. Firstly, data collection and preparation involve extensively gathering data from three types of graphs, resulting in a total of 15 sub-tasks in 29 datasets. These datasets are categorized into three types: conventional datasets for training and testing, counter-task datasets also used for training and testing to prevent model overfitting and enhance generalization, and out-of-distribution (OOD) datasets used solely for testing. Secondly, the three-stage curriculum learning fine-tuning framework, built upon a base model, includes the *KG Empowerment Stage*, which leverages KG datasets, the *EKG Enhancement Stage*, utilizing EKG datasets, and the *CKG Generalization Stage*, which incorporates CKG datasets along with counter-task datasets. Through these three stages of training, we obtain the micro, mid, and macro versions of GKG-LLM, respectively. Finally, GKG-LLM has undergone extensive testing and analysis on all three graph types across in-domain, OOD, and counter-task data, demonstrating the effectiveness and advancement of diverse instruction design strategies and the three-stage fine-tuning framework.

The contributions of this research are listed as follows:

- We propose an approach for building GKG using a three-stage curriculum learning fine-tuning framework, resulting in a GKG-LLM[1] that addresses task-specific differences and enables the unified construction of GKG.

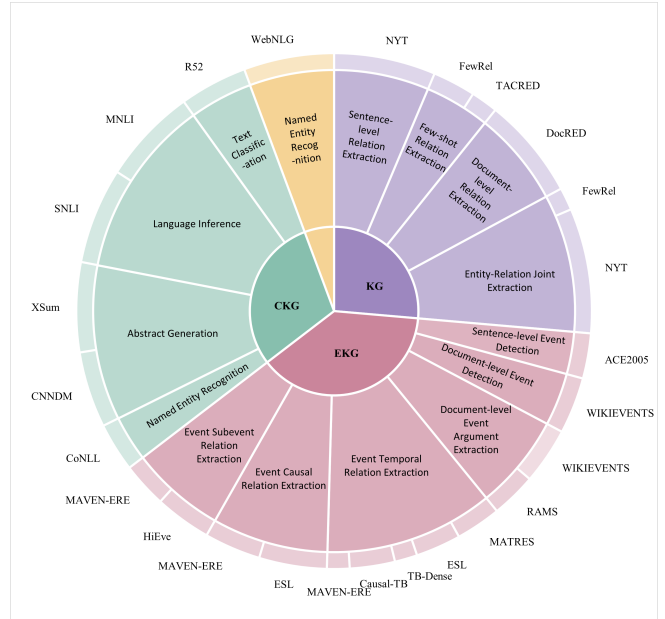- From a data perspective, this study is the first to collect

Figure 2: The illustration of the data distribution for all GKG sub-tasks.

and process sub-task datasets from three types of graphs in a comprehensive view, exploring their intrinsic connections in constructing GKG, as far as we know.

- Extensive experiments report that GKG-LLM achieves the effectiveness and advancement on three types of data and further analysis validates the superiority of our architecture.

## 2 Methodology

In this section, we first present the three-stage curriculum learning tuning framework in Section 2.1, then describe data collection and preparation in Section 2.2 and introduce our training strategy in Section 2.3.

The formal definition of GKG construction involves reformulating the various sub-tasks of KG, EKG, and CKG using a unified seq2seq format and structure. Then we solve it through three-stage fine-tuning LLMs, as shown in Figure 3. Specifically, the unified input is a task document or sentence, and the unified output consists of the elements or relations that form the GKG triples.

### 2.1 GKG-LLM

The overview of GKG-LLM is shown in Figure 3. It consists of three stages of tuning curriculum learning. Curriculum learning [Wang *et al.*, 2021] breaks down complex tasks into simpler ones and trains models in an increasing order of difficulty. This approach mimics the way humans learn by first mastering basic concepts before progressing to more complex knowledge.

From the previous theoretical analysis, we find that the three types of graphs have a progressive relationship. In a KG, entities and relations are represented as triples, which can be understood as event nodes in an EKG to some extent.
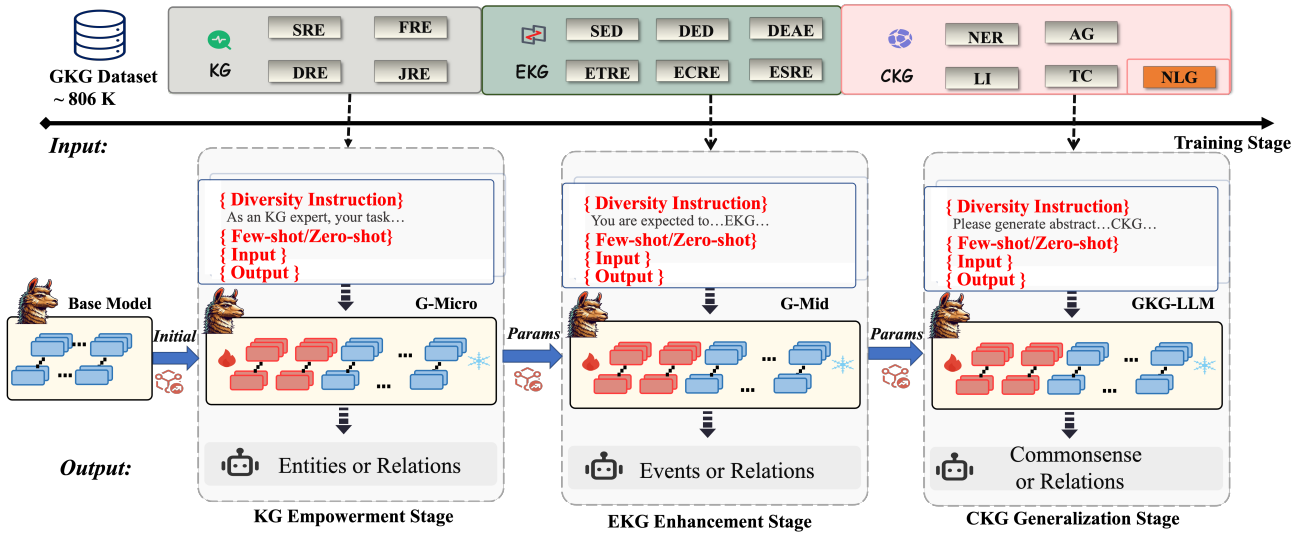
Figure 3: Three-stage curriculum learning tuning framework of GKG-LLM. The upper part represents the GKG dataset $\mathcal{D}_G$, consisting of the unified datasets. The lower part shows the three stages of GKG training: the *KG empowerment stage* using the KG datasets to build foundational skills, the *EKG enhancement stage* using the EKG datasets to enhance specific capabilities, and the *CKG generalization stage* using the CKG datasets and the counter task dataset to achieve generalization of the GKG-LLM capabilities. The thick arrows between the stages represent the delivery of model parameters from base model to each version of GKG-LLM.

EKG further explores the relationships between event nodes, while a CKG can be seen as a generalization of EKG, based on more universal commonsense knowledge.

Therefore, the tuning framework is divided into three stages following a curriculum learning approach: the *KG empowerment stage*, the *EKG enhancement stage*, and the *CKG generalization stage*. After the KG empowerment stage, we obtain the G-Micro model, which is expected to handle basic sub-tasks related to KG, such as handling various entity and relation extraction tasks. However, GKG nodes and relationships may include dynamic knowledge. Next, in the EKG enhancement stage, we utilize EKG-related sub-tasks datasets to further empower GKG-LLM on the basis of G-Micro, resulting in the G-Mid model, capable of handling sub-tasks involving dynamic knowledge. Furthermore, in the CKG generalization stage, we inject CKG-related sub-tasks and counter task data into the G-Mid model, generalizing the task handling capability of KG to broader scenarios, ultimately resulting in the GKG-LLM model.

**KG empowerment stage**    At this stage, we only inject the KG sub-task dataset into LLMs, and the training loss function is defined as cross-entropy loss:

$$\mathcal{L}_{\text{CE}} = -\sum_i p\left(y_i\right) \log p_\theta\left(\hat{y}_i \mid s_i; x_i\right), \qquad (1)$$

where $p_\theta$ represents the tunable LLM with parameters $\theta$, initialized from the base model. The instruction $s_i$ is concatenated with the input $x_i$ denotes the prompt format to LLMs. $\hat{y}_i$ is the predicted output, while $y_i$ represents the ground truth.

**EKG Enhancement Stage**    At this stage, we inject knowledge about dynamic nodes and relationships to enhance the model's capability. Specifically, we train the G-Micro model

from the first stage using the EKG sub-task dataset. This process expands the model's understanding of complex graphs, enabling it to handle dynamic nodes and relationships with temporal dependencies and causal features, improving its adaptability to changing data and laying a foundation for the subsequent stages. The loss function is the same as in the first stage.

**CKG Generalization Stage**    Real-world scenarios go beyond static knowledge and specific events, encompassing commonsense knowledge for a broader understanding. Therefore, at this stage, we train the G-Mid model from the second stage using the CKG sub-task dataset to enhance its generalization and applicability. This expands the model's commonsense knowledge, enabling it to excel in open-ended and complex reasoning tasks [Xu *et al.*, 2025]. The model becomes more practical and effective in real-world scenarios, ultimately resulting in the GKG-LLM.

This study conducts extensive testing and analysis on three types of data: In-domain, OOD and counter task data. Detailed implementation specifics is discussed in the following sections.

## 2.2 Data Collection and Preparation

As a comprehensive dataset encompassing the GKG construction tasks, it requires extensive datasets for each sub-task across the three types of graphs. Additionally, it is necessary to perform reasonable partitioning of the various datasets and format them to prepare for the unified GKG construction framework.

The overview of data distribution of all of GKG sub-tasks is shown as Figure 2. The GKG dataset is $\mathcal{D}_G = \mathcal{D}_{KG} \bigcup \mathcal{D}_{EKG} \bigcup \mathcal{D}_{CKG} \bigcup \mathcal{D}_{ct}$. Here, $\mathcal{D}_{KG}$ includes the sub-tasks of KG such as relation extraction and entity-relation joint extraction; For $\mathcal{D}_{EKG}$, sub-tasks include sentence-level

event detection, document-level event argument extraction, and event temporal relation extraction; And for $\mathcal{D}_{CKG}$, sub-tasks include summary generation and text inference. $\mathcal{D}_{ct}$ refers to a structure-to-text dataset, specifically the WebNLG task and dataset used for natural language generation, designed to serve as a counter-task for all GKG sub-tasks to prevent overfitting and enhance generalization without compromising the primary performance. Finally, we obtain $\mathcal{D}_G$ of $\sim$806K pieces for training and $\sim$140K pieces for testing. Details of each dataset are attached in Appendix A. The details of each sub-task are provided in Appendix F.

After data collection, we format each piece $i$ of the GKG dataset into a unified format, which includes $ID$, instruction $s_i$, few-shot $fs$ / zero-shot $zs$ , input $x_i$, and output $y_i$. Details of the data format and few-shot organization can be found in Appendix B.

### 2.3 Training Strategy

To effectively fine-tune our model on the unified dataset, we employ the LoRA+ [Hayou *et al.*, 2024] technique, an advanced version of Low-Rank Adaptation (LoRA), which has shown great promise in parameter-efficient fine-tuning (PEFT). LoRA+ adapts only a small subset of model parameters, reducing computational costs while maintaining high performance. By leveraging low-rank matrix approximations, LoRA+ allows us to efficiently update the model parameters without the need for extensive computational resources. Formally, LoRA+ modifies the weight matrix $W$ in the neural network as follows:

$$W' = W + \Delta W, \tag{2}$$

where $\Delta W = AB$, with $A \in \mathbb{R}^{d \times r}$ and $B \in \mathbb{R}^{r \times k}$. Here, $d$ is the dimension of the input, $k$ is the dimension of the output, and $r$ is the rank of the adaptation matrices, which is much smaller than both $d$ and $k$, making the adaptation parameter-efficient. To make better use of limited resources for training the model, the advancement of LoRA+ is reflected, as shown in Equation 3, in the use of different update hyperparameters $\eta_A$ and $\eta_B$ for the two low-rank matrices $A$ and $B$:

$$\begin{cases} A = A - \eta_A G_A \\ B = B - \eta_B G_B. \end{cases} \tag{3}$$

This approach accelerates convergence and effectively demonstrates the efficient and adaptive capabilities of GKG-LLM in handling GKG construction sub-tasks.

In summary, our training process harnesses the strengths of LoRA+ for efficient fine-tuning while experimenting with diverse data utilization strategies to optimize model performance for comprehensive GKG construction. This approach ensures that our model not only learns effectively from the data but also adapts seamlessly to various NLP tasks within GKG.

## 3 Experiments

In this section, we thoroughly evaluate the performance of GKG-LLM across three data settings, including in-sample data, counter-task data, and out-of-distribution data. The

baseline methods and evaluation metrics are presented in Section 3.1, while the main experimental results are presented in Sections 3.2. The stage generalization results are presented in Appendix C. Hyper-parameter settings are provided in Appendix E.

### 3.1 Baselines and Metrics

To perform a comprehensive evaluation, the final version of GKG-LLM is compared with two main categories of existing baselines: close-source baselines and open-source baselines.

For closed-source baselines, we access the model through the OpenAI API, specifically using the gpt-4-turbo-preview version[2], and the Anthropic API to access the Claude-3-Opus version[3] for evaluation. We also use the Google API to access the Gemini-1.5-Pro version[4] for evaluation.

For open-source baselines, we conduct experiments on two foundations: LlaMA-2-Chat[5] and LlaMA-3-Instruct[6]. The LlaMA-2-GKG is fine-tuned from Llama-2-Chat, while LlaMA-3-Instruct serves as the foundation for GKG-LLM and also acts as a baseline. This model is fine-tuned to fit a specific graph, serving as a strong baseline. Our integrated SFT method trains all datasets from the three types of graphs simultaneously.

Referencing the general evaluation metrics for each sub-task, for abstraction generation and structure-to-text tasks, the Rough-L metric is used, while all other tasks employ the F1 score as the evaluation metric.

### 3.2 Main Results

In this section, we thoroughly evaluate the performance of GKG-LLM on in-domain, OOD, and counter tasks. Specifically, as detailed in Table 1, we assess its performance across various sub-tasks in the three types of graphs. Compared to the baseline, the results demonstrate the effectiveness and practicality of GKG-LLM on the construction of all three graph types across in-domain, OOD, and counter-task data.

**KG Sub-task Datasets**  KG sub-task datasets focus on various types of relation extraction, including sentence-level relation extraction, few-shot relation extraction, and entity relation extraction, etc. Compared to the three closed-source LLMs, GKG-LLM achieves the best performance, with a minimum performance improvement of 12.04%. Additionally, when compared to a model tuned solely with KG sub-task datasets, GKG-LLM demonstrates a minimum performance gain of 7.6%. Across all baselines, GKG-LLM consistently achieves either the best or the second-best performance.

**EKG Sub-task Datasets**  EKG sub-task datasets primarily include event detection, event argument extraction, and event relation extraction. Compared to the three closed-source LLMs, GKG-LLM achieves the best performance, with a

---

| Graphs | Tasks | Datasets | GPT-4 | Claude-3 | Gemini-1.5 | LlaMA-2-GKG | LlaMA-3-Instruct | Single-SFT | Integrated-SFT | GKG-LLM |
|---|---|---|---|---|---|---|---|---|---|---|
| **KG** | SRE | NYT | 64.94 | 66.76 | 68.59 | 78.18 | 55.12 | 74.39 | _79.32_ | **80.63** |
| | FRE | FewRel | 26.28 | 27.45 | 30.20 | _89.45_ | 22.64 | 78.65 | 86.74 | **90.48** |
| | | TACRED | 18.85 | 20.23 | 22.43 | _86.71_ | 12.74 | 70.66 | 84.66 | **88.96** |
| | DRE | DOCRED | 38.84 | 36.28 | 42.63 | 83.18 | 34.63 | 74.53 | _83.61_ | **85.71** |
| | JE&RE | FewRel | 6.32 | 5.44 | 7.52 | **42.05** | 3.20 | 26.76 | 30.56 | _34.32_ |
| | | NYT | 6.22 | 5.85 | 8.36 | **53.33** | 0.0 | 40.16 | 48.66 | _52.27_ |
| **EKG** | SED | ACE2005 | 17.50 | 8.57 | 22.40 | 32.47 | 0.0 | 22.74 | _34.32_ | **80.63** |
| | DED | WIKIEVENTS | 16.54 | 9.14 | 14.87 | 24.87 | 18.62 | _29.59_ | 23.84 | **39.86** |
| | DEAE | WIKIEVENTS | 42.58 | 53.41 | 47.69 | _70.46_ | 41.76 | 63.38 | 69.30 | **75.22** |
| | | RAMS | 13.84 | 5.70 | 38.49 | 48.33 | 30.74 | _53.43_ | 52.09 | **63.62** |
| | | MATRES | 39.97 | 36.62 | 38.51 | _62.94_ | 22.79 | 37.91 | 44.26 | **71.51** |
| | ETRE | ESL | 64.24 | 47.65 | 42.18 | 68.96 | 21.67 | _74.06_ | 67.63 | **75.33** |
| | | TB-Dense | 43.73 | 36.58 | 42.43 | _52.89_ | 36.55 | 49.30 | 51.23 | **53.54** |
| | | Causal-TB | 6.67 | 8.01 | 8.74 | 42.79 | 16.43 | 37.35 | **49.83** | _45.26_ |
| | | MAVEN-ERE | 43.80 | 21.73 | 42.10 | 71.55 | 40.29 | 37.35 | _75.44_ | **81.95** |
| | | TCR* | 15.43 | 18.74 | _25.34_ | 24.88 | 24.71 | 20.68 | 22.09 | **26.45** |
| | ECRE | ESL | 28.57 | 19.26 | 55.21 | 75.33 | 26.33 | 62.92 | _78.74_ | **84.89** |
| | | MAVEN-ERE | 51.98 | 11.36 | 43.38 | 76.48 | 13.37 | 78.91 | _88.59_ | **90.18** |
| | | Causal-TB* | 39.67 | 41.23 | 43.44 | 33.94 | 30.02 | 48.41 | _48.80_ | **55.79** |
| | ESRE | HiEve | 38.81 | 30.92 | 48.83 | 55.60 | 48.61 | 57.64 | _58.01_ | **58.61** |
| | | MAVEN-ERE | 40.09 | 13.12 | 38.09 | _44.37_ | 33.49 | 39.11 | 37.30 | **48.49** |
| **CKG** | NER | CoNLL | 15.94 | 14.46 | 18.27 | _77.50_ | 15.60 | 64.74 | 70.53 | **82.30** |
| | AG† | CNNDM | 30 | 28 | 22 | _36_ | 18 | 35 | 35 | **45** |
| | | XSum | _33_ | 26 | 29 | 28 | 9 | 24 | 30 | **38** |
| | LI | SNLI | 51.26 | 47.56 | 60.38 | 69.51 | 44.50 | 87.09 | **89.35** | _89.03_ |
| | | MNLI | 81.80 | 39.33 | 48.80 | 58.97 | 53.70 | **86.78** | 84.62 | _86.35_ |
| | TC | R8* | **72.26** | 36.43 | 66.58 | 65.27 | 58.89 | 28.83 | 58.64 | _69.33_ |
| | | R52 | 82.18 | 83.75 | 80.63 | **94.16** | 29.68 | 89.02 | 88.81 | _90.34_ |
| _Counter_ | NLG† | WebNLG | 78 | 65 | 76 | _83_ | 15 | 80 | 80 | **85** |
| **Average Performance** | | | 38.25 | 29.81 | 39.07 | 59.70 | 26.83 | 52.97 | _60.41_ | **67.90** |

Table 1: Performance comparison across various datasets and tasks. The best result for each sub-task is highlighted in bold, while the second-best result is underlined. The OOD datasets are starred by *. † means the task is evaluated by metric Rough-L of percentage. The results for GPT-4, Claude-3, and Gemini-1.5 are obtained via their respective APIs. LlaMA-2-GKG, LlaMA-3-Instruct, Single-SFT, and Integrated-SFT are implemented by us. The GKG-LLM column represents the final model obtained after three-stage tuning.

minimum improvement of 9.88%. An interesting observation is that the Integrated SFT model achieves the second-best performance in half of the tasks; however, GKG-LLM still consistently performs either the best or the second-best overall. Another interesting point is that in the OOD datasets, specifically the TCR dataset for the ETRE sub-task and the Causal-TB dataset for the ECRE sub-task, GKG-LLM outperforms the second-best baseline by 1.11% and 6.99%, respectively, demonstrating its strong generalization capability on OOD data.

**CKG Sub-task Datasets** For the CKG sub-task dataset, the focus is closer to common-sense nodes and relations reasoning, involving tasks such as abstract generation and language inference. For the R8 dataset in the Text Classification sub-task, which serves as an OOD dataset, GPT-4 achieves the best performance, attributed to its exceptional capabilities in language understanding. Even so, GKG-LLM still achieves the second-best performance. Since CKG closely resembles

real-world commonsense scenarios, both LlaMA-2-GKG and Single-SFT also demonstrates strong results. However, overall, GKG-LLM consistently maintains either the best or the second-best performance.

GKG-LLM achieves the best performance on the WebNLG dataset for the Natural Language Generation (NLG) task, surpassing the strongest baseline by 2%, further highlighting its strong structure-to-text capabilities. It consistently performs at the best or second-best level across all GKG sub-tasks, with an average improvement of 7.49% over the strongest baseline. Additionally, its strong performance on OOD data demonstrates its ability to generalize effectively to unseen data distributions, with ablation studies and OOD analysis detailed in Section 4.

### 3.3 Exploration of Three Stages

As discussed in Section 1, a triple in a KG can, to some extent, be considered as a node in an EKG, while the triples in EKG and CKG are linked through the relationship between

the concrete and the abstract. Theoretically, there exists a progressive relationship among these three types of graphs, which serves as the theoretical basis for our three-stage fine-tuning framework. Therefore, this subsection will explore the performance of the three types of graphs under different fine-tuning sequences, as well as the performance of the intermediate versions of our three-stage fine-tuning framework on the sub-tasks of the three types of graphs.
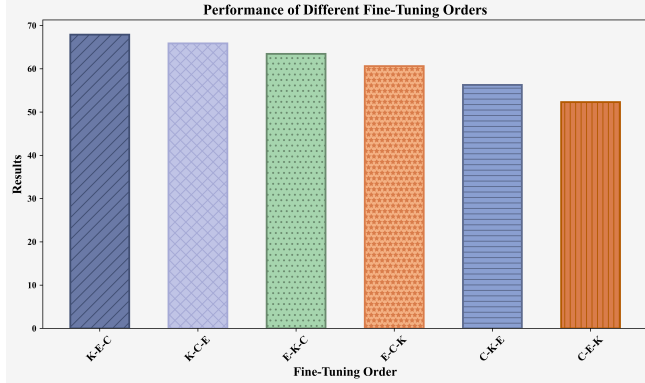


Figure 4: Results of different fine-tuning orders. "K-E-C" means the fine-tuning order is KG, EKG and CKG. The following sets of experiments are similar to this one.

As shown in Figure 4, the three types of graphs show varying performance in terms of average performance across all tasks under different fine-tuning sequences. The "K-E-C" sequence adopted in this study demonstrates the best performance, further confirming the theoretical correctness and experimental effectiveness of our three-stage fine-tuning sequence.
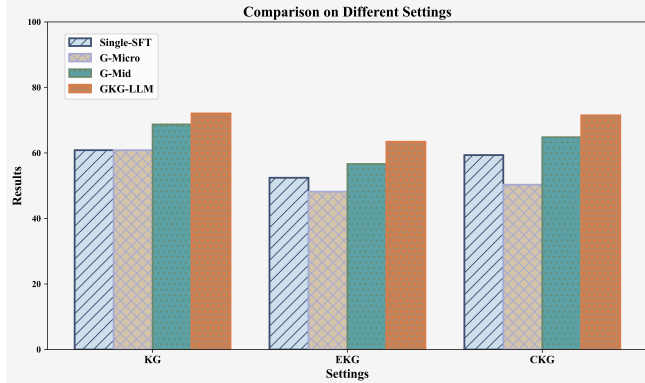


Figure 5: Fine-tuning with a single type of graph and performance of different intermediate version in the GKG-LLM.

Figure 5 presents the performance of the single SFT model and the three-stage models across the KG, EKG, and CKG sub-tasks. In each sub-task, the results improve as the fine-tuning progresses through the three stages. Compared to single-SFT, our GKG-LLM framework demonstrates better performance, validating the practicality of the three-stage fine-tuning approach.

# 4 Analysis

In this section, we introduce the ablation study in Section 4.1 and provide a comprehensive analysis and explanation of the OOD data in Section 4.2. An analysis of data scaling in training is introduced in Section 4.3. The evaluation of the optimal model under various hyper-parameter settings is presented in Appendix D.

| Variation | KG | EKG | CKG | Avg. |
|---|---|---|---|---|
| $\mathcal{P}$ | 72.06 | 63.42 | 71.48 | 67.90 |
| $\mathcal{P}_{si}$ | 68.46 | 59.34 | 69.10 | 64.33 |
| $\Delta$ | (-3.60) | (-4.08) | (-2.38) | (-3.57) |
| $\mathcal{P}_{zs}$ | 65.17 | 55.09 | 66.05 | 60.06 |
| $\Delta$ | (-6.89) | (-8.33) | (-5.43) | (-7.84) |
| $\mathcal{P}_{si+zs}$ | 62.44 | 52.26 | 64.66 | 58.15 |
| $\Delta$ | (-9.62) | (-11.16) | (-6.82) | (-9.75) |

Table 2: Performance comparison of different prompt strategies on the evaluation metrics. $\mathcal{P}$ denotes full prompts, $\mathcal{P}_{si}$ refers to a single instruction regardless of diversity, $\mathcal{P}_{zs}$ represents zero-shot only, and $\mathcal{P}_{si+zs}$ combines single instruction with zero-shot prompting.

## 4.1 Ablation Studies

In this section, we present the ablation study for three different prompt strategies: (1) using only a single instruction to construct the prompt format, (2) using only zero-shot prompts without employing any few-shot examples, and (3) removing both strategies simultaneously. We compare the performance across three types of graphs and the overall dataset, with the comparison results shown in Table 2. Examples of different types of prompts can be found in the respective sections of Appendix B.

The results show that removing the diversity of instructions causes a noticeable performance drop, as diverse instructions better reflect real-world scenarios where different questioners have unique styles, requiring the model to adapt to various instruction formats. Removing the few-shot learning strategy lead to an even greater performance degradation, as LLMs lost their ability to perform in-context learning and relies only on inherent capabilities, affecting their ability to generate the corresponding elements or relationships. The most performance drop occurs when both strategies are removed, highlighting that the advantages of these strategies are cumulative, further validating the superiority and effectiveness of our data construction strategy.

## 4.2 OOD Analysis

This section specifically discusses the performance of GKG-LLM on OOD datasets. As introduced in Section 2.1, our data is divided into three parts, with the OOD portion deliberately excluded during the initial training design, meaning that GKG-LLM has never encountered these types of data before. Therefore, the performance on this part serves as an indicator of our model's generalization ability from the perspective of OOD data.

As shown in Figure 7, overall, our method achieves the best performance, reaching 50.52%, which is 5.40% higher

than the second-best model, Gemini-1.5-pro. Despite the fact that these data points were entirely unfamiliar to both closed-source LLMs and our tuned open-source LLMs, our model still demonstrates strong robustness and effectiveness.

## 4.3 Analysis on Different Data Scaling

This section explores the impact of different data scales on model performance. The model is trained using 10%, 20%, 40%, 60%, 80%, and 100% of the data, sampled from the three types of graph sub-tasks separately. The results show that as the data proportion increases, model performance improves progressively, with performance being limited at 10%, improving at 20% and 40%, and continuing to enhance at 60% and 80%, reaching near-optimal performance at 100%.
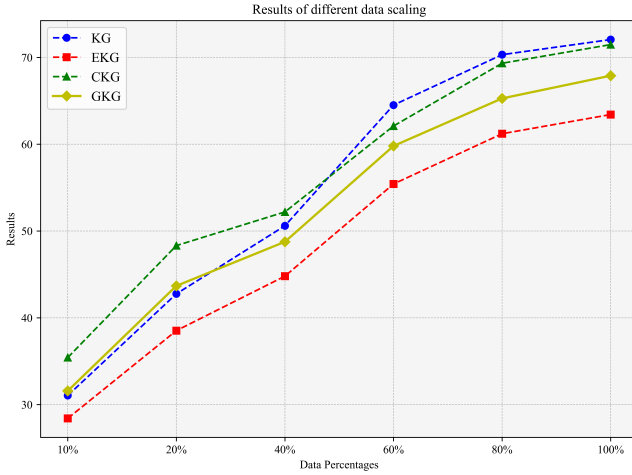


Figure 6: Results of training with different proportions of complete data.

Figure 6 shows that as the data volume increases, the model's average scores across all tasks gradually improve. Notably, the average scores for the three types of graph sub-tasks follow similar trends, with diminishing performance gains beyond 80% data usage, indicating a saturation point where the additional data brings marginal benefits.

## 5 Related Works

This section introduces two types of related work. Section 5.1 covers three typical tasks within GKG sub-tasks, while Section 5.2 discusses research related to LLMs.

### 5.1 GKG Sub-tasks

In this section, we introduce a representative task for each of the three types of graphs: the entity-relation joint extraction task in the KGs, the document-level event argument extraction task in the EKGs, and the abstract generation task in the CKGs.

**Entity-relation joint extraction** task has been a focus in the domain of knowledge graph construction, as it aims to simultaneously extract entities and their relationships from unstructured text. Current state-of-the-art methods leverage
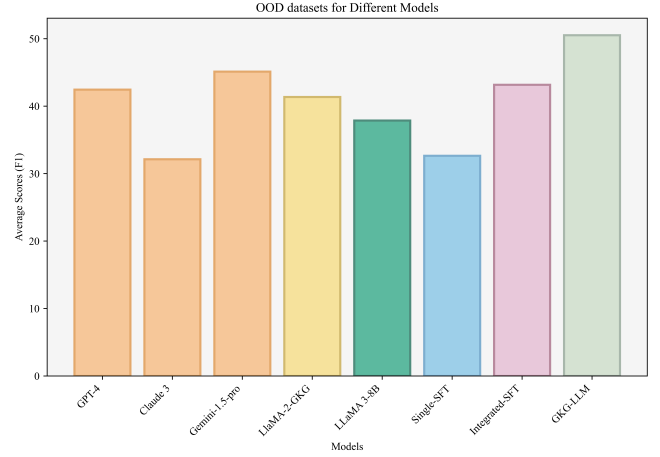


Figure 7: The average performance on OOD datasets, consisting TCR, Causal-TB and R8 datasets.

transformer architecture to model interactions between entities within sentences or documents, which provides further performance gains [Sui *et al.*, 2023]. **Document-level event argument extraction** aims to extract the arguments of events from long texts to better understand complex event relations and event chains. Pre-trained models such as BERT have been widely employed in event extraction tasks. By combining pre-trained knowledge with task-specific fine-tuning, these models have proven effective in understanding complex contexts [Zhang *et al.*, 2024]. **Abstract generation** particularly with the rise of pre-trained transformer-based models. A recent state-of-the-art approach by [Gao *et al.*, 2023] utilizes a combination of pre-trained language models and reinforcement learning to enhance the quality of generated abstracts.

### 5.2 Large Language Models

With the emergence of closed-source and open-source LLMs represented by GPT4 [Achiam *et al.*, 2023] and LlaMA-3 [Dubey *et al.*, 2024], respectively, a large amount of research has focused on these models. This section introduces some of the work based on close-source and open-source LLMs.

Research based on closed-source LLMs typically involves evaluating these large models [Gandhi *et al.*, 2024] and integrating them with traditional tasks. For example, such studies may focus on enhancing certain aspects of conventional natural language tasks [Zheng *et al.*, 2023] or providing new perspectives for text analysis [Savelka *et al.*, 2023]. The study by [Xu *et al.*, 2024] using LlaMA-2 as the foundation, explores the possibility of a unified approach to symbol-centric tasks through full fine-tuning and extend this approach to generalize to natural language-centric tasks. A survey by [Zhang *et al.*, 2023] introduce various paradigms of instruction fine-tuning for LLMs, providing a comprehensive overview of its advantages, limitations, and implementation methods.

However, up to now, no study has integrated the broad task of GKG construction. This research unifies such tasks from both the task and data perspectives by fine-tuning open-source LLMs.

# 6 Conclusion

This study proposes a new task for building GKG. It represents the first collection approached from the unified perspective in terms of data, and the first unified construction of three types of graphs from the task perspective. This task addresses two issues: obstacles arising from differences between tasks, and the neglect of intrinsic connections among different types of graphs. To address these challenges, we propose a three-stage curriculum learning framework that iteratively injects sub-task knowledge from KG, EKG, and CKG into GKG-LLM, aiming for broad and outstanding performance in GKG construction. Extensive experiments demonstrate the effectiveness and robustness of the GKG-LLM approach. The models and data from this study will be fully released upon acceptance of the paper. In the future, we will expand the application of GKG-LLM into a broader range of scenarios, such as intelligent healthcare [He *et al.*, 2025; Lin *et al.*, 2025b], to enhance its utility and impact.

## References

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

[Alt *et al.*, 2020] Christoph Alt, Aleksandra Gabryszak, and Leonhard Hennig. Tacred revisited: A thorough evaluation of the tacred relation extraction task. *arXiv preprint arXiv:2004.14855*, 2020.

[Camburu *et al.*, 2018] Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31, 2018.

[Chan *et al.*, 2024] Chunkit Chan, Cheng Jiayang, Weiqi Wang, Yuxin Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. Exploring the potential of chatgpt on sentence level relations: A focus on temporal, causal, and discourse relations. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 684–721, 2024.

[Chen *et al.*, 2021] Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. Dialogsum: A real-life scenario dialogue summarization dataset. *arXiv preprint arXiv:2105.06762*, 2021.

[Dubey *et al.*, 2024] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[Ebner *et al.*, 2020] Seth Ebner, Patrick Xia, Ryan Culkin, Kyle Rawlins, and Benjamin Van Durme. Multi-sentence argument linking. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8057–8077, 2020.

[Gandhi *et al.*, 2024] Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. Understanding social reasoning in language models with language models. *Advances in Neural Information Processing Systems*, 36, 2024.

[Gao *et al.*, 2023] Catherine A Gao, Frederick M Howard, Nikolay S Markov, Emma C Dyer, Siddhi Ramesh, Yuan Luo, and Alexander T Pearson. Comparing scientific abstracts generated by chatgpt to real abstracts with detectors and blinded human reviewers. *NPJ Digital Medicine*, 6(1):75, 2023.

[Gardent *et al.*, 2017] Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. The webnlg challenge: Generating text from rdf data. In *10th International Conference on Natural Language Generation*, pages 124–133. ACL Anthology, 2017.

[Ge and Moh, 2017] Lihao Ge and Teng-Sheng Moh. Improving text classification with word embedding. In *2017 IEEE International Conference on Big Data (Big Data)*, pages 1796–1805. IEEE, 2017.

[Glavaš *et al.*, 2014] Goran Glavaš, Jan Šnajder, Parisa Kordjamshidi, and Marie-Francine Moens. Hieve: A corpus for extracting event hierarchies from news stories. 2014.

[Grishman *et al.*, 2005] Ralph Grishman, David Westbrook, and Adam Meyers. Nyu's english ace 2005 system description. *Ace*, 5(2), 2005.

[Gubelmann *et al.*, 2024] Reto Gubelmann, Ioannis Katis, Christina Niklaus, and Siegfried Handschuh. Capturing the varieties of natural language inference: A systematic survey of existing datasets and two novel benchmarks. *Journal of Logic, Language and Information*, 33(1):21–48, 2024.

[Han *et al.*, 2018] Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation. *arXiv preprint arXiv:1810.10147*, 2018.

[Han *et al.*, 2019] Rujun Han, I Hsu, Mu Yang, Aram Galstyan, Ralph Weischedel, Nanyun Peng, et al. Deep structured neural network for event temporal relation extraction. *arXiv preprint arXiv:1909.10094*, 2019.

[Hasan *et al.*, 2021] Tahmid Hasan, Abhik Bhattacharjee, Md Saiful Islam, Kazi Samin, Yuan-Fang Li, Yong-Bin Kang, M Sohel Rahman, and Rifat Shahriyar. Xl-sum: Large-scale multilingual abstractive summarization for 44 languages. *arXiv preprint arXiv:2106.13822*, 2021.

[Hayou *et al.*, 2024] Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. *arXiv preprint arXiv:2402.12354*, 2024.

[He *et al.*, 2022] Yong He, Cheng Wang, Shun Zhang, Nan Li, Zhaorong Li, and Zhenyu Zeng. Kg-mtt-bert: Knowledge graph enhanced bert for multi-type medical text classification. *arXiv preprint arXiv:2210.03970*, 2022.

[He *et al.*, 2025] Kai He, Rui Mao, Qika Lin, Yucheng Ruan, Xiang Lan, Mengling Feng, and Erik Cambria. A survey

of large language models for healthcare: from data, technology, and applications to accountability and ethics. *Information Fusion*, 118:102963, 2025.

[Hettiarachchi *et al.*, 2023] Hansi Hettiarachchi, Mariam Adedoyin-Olowe, Jagdev Bhogal, and Mohamed Medhat Gaber. Ttl: transformer-based two-phase transfer learning for cross-lingual news event detection. *International Journal of Machine Learning and Cybernetics*, 2023.

[Hu *et al.*, 2020] Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra Kübler, and Lawrence S Moss. Ocnli: Original chinese natural language inference. *arXiv preprint arXiv:2010.05444*, 2020.

[Huang *et al.*, 2020] Kung-Hsiang Huang, Mu Yang, and Nanyun Peng. Biomedical event extraction with hierarchical knowledge graphs. *arXiv preprint arXiv:2009.09335*, 2020.

[Krause *et al.*, 2022] Franz Krause, Tobias Weller, and Heiko Paulheim. On a generalized framework for time-aware knowledge graphs. In *Towards a Knowledge-Aware AI*, pages 69–74. IOS Press, 2022.

[Lai *et al.*, 2023] Vivian Lai, Chacha Chen, Alison Smith-Renner, Q Vera Liao, and Chenhao Tan. Towards a science of human-ai decision making: An overview of design space in empirical human-subject studies. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1369–1385, 2023.

[Li *et al.*, 2021] Sha Li, Heng Ji, and Jiawei Han. Document-level event argument extraction by conditional generation. *arXiv preprint arXiv:2104.05919*, 2021.

[Lin *et al.*, 2023] Qika Lin, Jun Liu, Rui Mao, Fangzhi Xu, and Erik Cambria. TECHS: temporal logical graph networks for explainable extrapolation reasoning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1281–1293, 2023.

[Lin *et al.*, 2025a] Qika Lin, Tianzhe Zhao, Kai He, Zhen Peng, Fangzhi Xu, Ling Huang, Jingying Ma, and Mengling Feng. Self-supervised quantized representation for seamlessly integrating knowledge graphs with large language models. *CoRR*, abs/2501.18119, 2025.

[Lin *et al.*, 2025b] Qika Lin, Yifan Zhu, Xin Mei, Ling Huang, Jingying Ma, Kai He, Zhen Peng, Erik Cambria, and Mengling Feng. Has multimodal learning delivered universal intelligence in healthcare? A comprehensive survey. *Information Fusion*, 116:102795, 2025.

[Ma *et al.*, 2023] Youmi Ma, An Wang, and Naoaki Okazaki. Dreeam: Guiding attention with evidence for improving document-level relation extraction. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1971–1983, 2023.

[Mirza and Tonelli, 2016] Paramita Mirza and Sara Tonelli. Catena: Causal and temporal relation extraction from natural language texts. In *The 26th international conference on computational linguistics*, pages 64–75. ACL, 2016.

[Ning *et al.*, 2019] Qiang Ning, Sanjay Subramanian, and Dan Roth. An improved neural baseline for temporal relation extraction. *arXiv preprint arXiv:1909.00429*, 2019.

[Paulus, 2017] R Paulus. A deep reinforced model for abstractive summarization. *arXiv preprint arXiv:1705.04304*, 2017.

[Peng *et al.*, 2023] Ciyuan Peng, Feng Xia, Mehdi Naseriparsa, and Francesco Osborne. Knowledge graphs: Opportunities and challenges. *Artificial Intelligence Review*, 56(11):13071–13102, 2023.

[Pimenov *et al.*, 2023] Danil Yu Pimenov, Andres Bustillo, Szymon Wojciechowski, Vishal S Sharma, Munish K Gupta, and Mustafa Kuntoğlu. Artificial intelligence systems for tool condition monitoring in machining: Analysis and critical review. *Journal of Intelligent Manufacturing*, 34(5):2079–2121, 2023.

[Sang and De Meulder, 2003] Erik F Sang and Fien De Meulder. Introduction to the conll-2003 shared task: Language-independent named entity recognition. *arXiv preprint cs/0306050*, 2003.

[Savelka *et al.*, 2023] Jaromir Savelka, Kevin D Ashley, Morgan A Gray, Hannes Westermann, and Huihui Xu. Can gpt-4 support analysis of textual data in tasks requiring highly specialized domain expertise? *arXiv preprint arXiv:2306.13906*, 2023.

[Sui *et al.*, 2023] Dianbo Sui, Xiangrong Zeng, Yubo Chen, Kang Liu, and Jun Zhao. Joint entity and relation extraction with set prediction networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[Wadhwa *et al.*, 2023] Somin Wadhwa, Silvio Amir, and Byron C Wallace. Revisiting relation extraction in the era of large language models. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2023, page 15566. NIH Public Access, 2023.

[Wang *et al.*, 2021] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):4555–4576, 2021.

[Wang *et al.*, 2022] Xiaozhi Wang, Yulin Chen, Ning Ding, Hao Peng, Zimu Wang, Yankai Lin, Xu Han, Lei Hou, Juanzi Li, Zhiyuan Liu, et al. Maven-ere: A unified large-scale dataset for event coreference, temporal, causal, and subevent relation extraction. *arXiv preprint arXiv:2211.07342*, 2022.

[Xu *et al.*, 2024] Fangzhi Xu, Zhiyong Wu, Qiushi Sun, Siyu Ren, Fei Yuan, Shuai Yuan, Qika Lin, Yu Qiao, and Jun Liu. Symbol-llm: Towards foundational symbol-centric interface for large language models. In *ACL*, 2024.

[Xu *et al.*, 2025] Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. Are large language models really good logical reasoners? a comprehensive evaluation and beyond. *IEEE Transactions on Knowledge and Data Engineering*, 2025.

[Yamada and Shindo, 2019] Ikuya Yamada and Hiroyuki Shindo. Neural attentive bag-of-entities model for text classification. *arXiv preprint arXiv:1909.01259*, 2019.

[Yao *et al.*, 2019] Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. Docred: A large-scale document-level relation extraction dataset. *arXiv preprint arXiv:1906.06127*, 2019.

[Zhang *et al.*, 2023] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.

[Zhang *et al.*, 2024] Jian Zhang, Changlin Yang, Haiping Zhu, Qika Lin, Fangzhi Xu, and Jun Liu. A semantic mention graph augmented model for document-level event argument extraction. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 1577–1587, 2024.

[Zheng *et al.*, 2023] Mingkai Zheng, Xiu Su, Shan You, Fei Wang, Chen Qian, Chang Xu, and Samuel Albanie. Can gpt-4 perform neural architecture search? *arXiv preprint arXiv:2304.10970*, 2023.

# A Details of Data Collection

This section provides detailed information on all datasets of ∼806K pieces for training and ∼140K pieces for testing, including an overall introduction in Section A.1, and the categorization of datasets into three types in Section A.2.

## A.1 General Introduction

As shown in Table 3, we have collected, to the best of our ability, three types of different graph construction sub-task datasets for the GKG Dataset, along with an additional counter task (NLG task) dataset, resulting in a total of 15 sub-tasks and 29 datasets. To ensure data balance and reasonable distribution, we sample and partition some of the datasets. These sampling and partitioning processes are clearly indicated in Table 3 under the "Sampled?" field, allowing readers to better understand the data handling approach.

In the KG sub-task dataset, the focus is primarily on various types of relation extraction, including sentence-level relation extraction, few-shot relation extraction, and entity relation extraction, etc. This is because nodes in the KG sub-task are entities, and an important sub-task is to extract relationships between these entities. Furthermore, the EKG sub-task dataset primarily includes event detection, event argument extraction, and event relation extraction, as the event nodes are more complex, containing trigger words and various arguments. For the CKG sub-task dataset, the focus is closer to common-sense nodes and relations reasoning, involving tasks such as abstract generation and language inference.
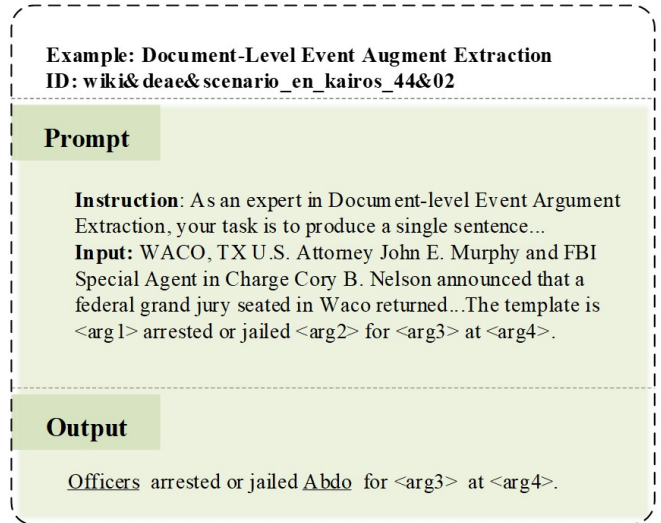


Figure 8: An example from the WIKEVENTS dataset. It consists of five fields $ID$, instruction $s_i$, few-shot $fs$ / zero-shot $zs$ , input $x_i$, and output $y_i$.

## A.2 Three Categorizations

The GKG Dataset is divided into three types: **in-domain data**, **counter task data**, and **OOD data**. The OOD data is separately indicated in Table 3 and is used only during the testing phase, not during training, to evaluate the model's performance on OOD data. The counter task is included to prevent overfitting and to enhance the generalizability of GKG-LLM.

Specifically, in-domain data consists of various GKG sub-tasks, combined with the counter task dataset (WebNLG) to form the training set. Using a curriculum learning fine-tuning framework, we obtained the final version of GKG-LLM. After testing on all in-domain datasets and the counter task dataset, we proceeded to test on three OOD datasets—TCR, Causal-TB, and R8—to validate the model's superior performance.

## B Data Format

To bridge the gap between the dataset's data format and the instruction-tuning format, we formatted all the data. Specifically, each data entry consists of five fields– $ID$, instruction $s_i$, few-shot $fs$ / zero-shot $zs$ , input $x_i$, and output $y_i$. as shown in Figure 8, this example is from the WIKIEVENTS dataset. $ID$ represents the unique identifier of each data entry, which includes the task name, dataset name, and specific data entry. The instruction $s_i$ provides a formal definition of each sub-task and is passed to the base model to help it understand the task's intent. few-shot $fs$ / zero-shot $zs$ field indicates whether a few-shot example is included in the prompt; in particular, for zero-shot, this field can be omitted. The input $x_i$ represents the specific input data, while the output $y_i$ represents the corresponding output.

To more comprehensively simulate real-world scenarios, we utilize GPT-4 to generate ten **diverse instructions**, which are then randomly assigned to the instruction field of each

| Graphs | Tasks | Datasets | # Train | # Test | sampled? | *held-out?* | Original Source |
|---|---|---|---|---|---|---|---|
| **KG** | SRE | NYT | 96,229 | 8,110 | | | [Paulus, 2017] |
| | FRE | FewRel | 56,576 | 11,775 | | | [Han *et al.*, 2018] |
| | | TACRED | 18,448 | 3,325 | | | [Alt *et al.*, 2020] |
| | DRE | DOCRED | 61,380 | 6,137 | ✓ | | [Yao *et al.*, 2019] |
| | JE&RE | FewRel | 28,288 | 11,775 | ✓ | | |
| | | NYT | 48,114 | 8,110 | ✓ | | |
| **EKG** | SED | ACE2005 | 3,681 | 409 | | | [Grishman *et al.*, 2005] |
| | DED | WIKIEVENTS | 3,586 | 365 | | | [Li *et al.*, 2021] |
| | DEAE | WIKIEVENTS | 3,586 | 365 | | | |
| | | RAMS | 7,339 | 761 | | | [Ebner *et al.*, 2020] |
| | | MATRES | 12,216 | 1,361 | | | [Ning *et al.*, 2019] |
| | | ESL | 7,652 | 852 | | | |
| | ETRE | TB-Dense | 9,257 | 2,639 | | | [Han *et al.*, 2019] |
| | | Causal-TB | 5,427 | 603 | | | [Mirza and Tonelli, 2016] |
| | | MAVEN-ERE | 80,000 | 5,000 | ✓ | | [Wang *et al.*, 2022] |
| | | TCR | | 3,515 | | ✓ | [Han *et al.*, 2019] |
| | | ESL | 3,196 | 356 | | | |
| | ECRE | MAVEN-ERE | 63,980 | 7,330 | ✓ | | |
| | | Causal-TB | | 318 | | ✓ | |
| | ESRE | HiEve | 12,107 | 1,348 | | | [Glavaš *et al.*, 2014] |
| | | MAVEN-ERE | 31,365 | 4,244 | | | |
| **CKG** | NER | CoNLL | 17,293 | 3,454 | | | [Sang and De Meulder, 2003] |
| | AG | CNNDM | 51,684 | 11,490 | ✓ | | [Chen *et al.*, 2021] |
| | | XSum | 50,666 | 11,334 | ✓ | | [Hasan *et al.*, 2021] |
| | LI | SNLI | 50,000 | 10,000 | ✓ | | [Camburu *et al.*, 2018] |
| | | MNLI | 50,000 | 10,000 | ✓ | | [Hu *et al.*, 2020] |
| | TC | R8 | | 7,674 | | ✓ | [Yamada and Shindo, 2019] |
| | | R52 | 7,816 | 1,284 | ✓ | | [Ge and Moh, 2017] |
| *Counter* | NLG | WebNLG | 26,302 | 6,513 | | | [Gardent *et al.*, 2017] |

Table 3: Detailed illustrations of 15 sub-task types across 29 datasets, categorized within three types of graphs, along with a counter dataset—WebNLG. # Train and # Test represent the number of training and testing samples, respectively. *Sampled?* indicates whether the dataset is sampled from the original to achieve data balancing. *Held-out?* specifies whether the dataset is used during the training phase. *Original Source* refers to the citation of the original paper.

data entry. This approach aims to enhance the model's ability to understand and handle a variety of task instructions, thereby increasing its flexibility and adaptability for real-world multitasking needs. By diversifying the instructions, we aim to train the model to better respond to different directives, similar to a practical deployment setting. Additionally, for 10% of the data pieces, we randomly added a **few-shot example** to help the base model understand the task structure more effectively. The majority of the data entries, however, remained in a zero-shot setting, ensuring that the model could learn general patterns of GKG construction tasks without extensive direct guidance. By balancing few-shot and zero-shot learning, we aim to improve the model's generalization capabilities across a range of GKG-related tasks.

## C Stage Generalization

In this section, we examine the effect of the three-stage training strategy on subsequent data exploration stages. Specifically, we test G-Micro, trained only on KG-related sub-task datasets, on EKG and CKG sub-task datasets, and G-Mid on the CKG sub-task dataset. The results are shown in Figure 9.
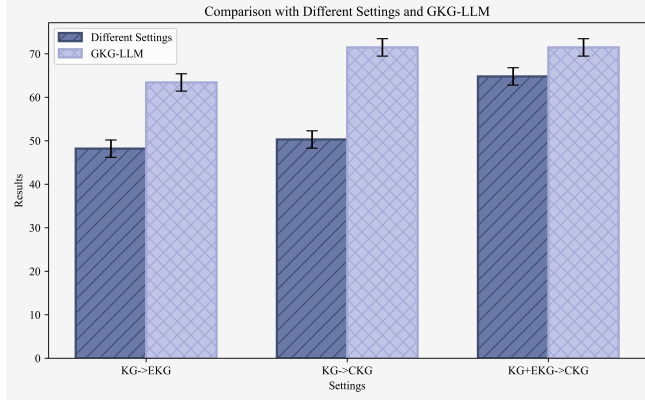


Figure 9: Comparison of Results by different settings and GKG-LLM.

The experimental results show that, despite some trade-offs in the exploratory experiments, the three-stage curriculum learning approach achieves superior performance. This demonstrates: **1).** earlier GKG-LLM versions influence subsequent tasks, indicating task correlation; **2).** the unified approach to the three types of graphs in GKG is valuable and meaningful, reflecting their progressive relationship within a unified framework.

## D Exploration of LoRA+ Hyperparameter Values

As described in Section 2.3, we adopt the LoRA+ training strategy, where the low-rank matrices $A$ and $B$ have different rates of change, meaning they each have distinct hyperparameters $\eta_A$ and $\eta_B$.

In this section, we explore the effects of different combinations of the hyperparameters $\eta_A$ and $\eta_B$ on the model's

performance. The experimental results are illustrated in Figure 10, the vertical axis represents $B$, which is expressed as a multiple of $\eta_A$. The model's performance is highly sensitive to changes in $\eta_A$ and $\eta_B$. The highest performance score of 67.90% was achieved with $\eta_A = 4 \times 10^{-4}$ and $\eta_B = 4 \times 10^{-3}$. This suggests that higher learning rates for $\eta_A$ combined with moderate values of $\eta_B$ are beneficial for fine-tuning. Conversely, the lowest performance scores were observed with the smallest value of $\eta_A = 5 \times 10^{-5}$, regardless of the value of $\eta_B$. This indicates that too low a learning rate for the adaptation matrices may not be sufficient for effective fine-tuning. Increasing $\eta_B$ tends to enhance performance up to a certain point, after which the performance gains stabilize or diminish. For example, $\eta_A = 2 \times 10^{-4}$ with $\eta_B = 8 \times 10^{-3}$ shows a obvious score, but further increasing $\eta_B$ does not yield substantial improvements.
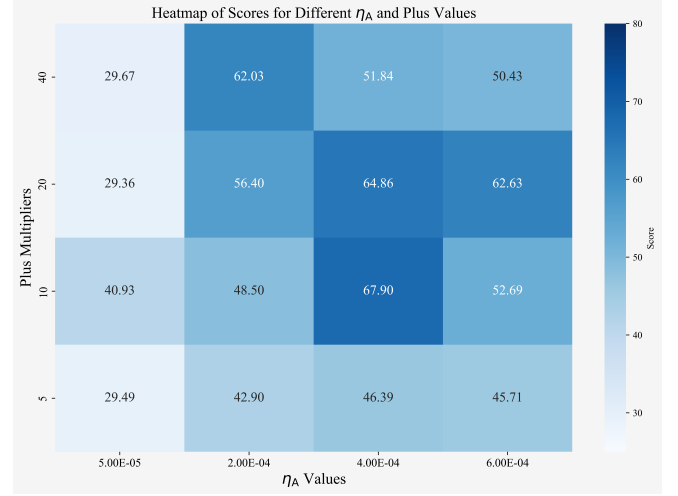


Figure 10: Heatmap of Scores for Different $\eta_A$ and $\eta_B$ Values for our training strategy.

These findings highlight the importance of carefully tuning the hyperparameters $\eta_A$ and $\eta_B$ in the LoRA+ framework to achieve optimal model performance. The insights gained from this exploration can guide future experiments and the development of more effective fine-tuning strategies for LLMs. In summary, the exploration of LoRA+ hyperparameters reveals that selecting the appropriate values for $\eta_A$ and $\eta_B$ is crucial for maximizing model performance. This study provides a foundational understanding that can be leveraged to further enhance the efficiency and effectiveness of fine-tuning LLMs using low-rank adaptation techniques.

## E Hyper-parameters

In the implementation, we leverage the LoRA+ technique to fine-tune models using four A800 (80GB) GPUs, with a maximum sequence length of 4,096. The fine-tuning process is optimized with FlashAttention2, while the AdamW optimizer is employed with a learning rate of 5e-5 across three curriculum learning stages, each controlled by a linear learning rate scheduler. We use one epoch per stage to complete the tuning process.

During the KG empowerment stage, model weights are initialized from LLaMA-3-Instruct, resulting in the tuned model named G-Micro. In the EKG enhancement stage, G-Micro serves as the starting point, producing G-Mid. Similarly, in the CKG generalization stage, we initialize from G-Mid and ultimately obtain GKG-LLM. Inference process is conduct on a single A800 (80GB) GPU using greedy search.

# F  Sub-tasks Introduction

The GKG dataset is composed of three types of sub-task datasets: KG , EKG and CKG. The data is categorized into three types: In-domain data, OOD data, and counter-task data. The specific descriptions of these tasks are as follows.

## F.1  KG

**SRE (Sentence-level Relation Extraction)**  For the SRE task, we utilize the NYT dataset. This task focuses on identifying the entities mentioned in a complex news sentence and, based on entity recognition, detecting and labeling the relationships between the entities. This task plays a critical role in the process of transforming unstructured textual data into structured knowledge.

**FRE (Few-shot Relation Extraction)**  Due to the issue of insufficient labeled corpora in many domains and the high cost of manual annotation, the FRE task aims to train a model using a small amount of labeled sample data, enabling the model to learn the characteristic information of entities that form relationships. During the testing phase, the model is asked to identify previously unseen relationship types from new datasets. In our work, we utilize the FewRel and TA-CRED datasets for both training and testing.

**DRE (Document-level Relation Extraction)**  Compared to SRE, the DRE task is more challenging, as it requires the model not only to identify relations within a single sentence but also to understand the context and possess the ability to recognize relations across sentences and even across paragraphs. In this paper, we conduct experiments using the DocRED dataset. The input is a long text document containing multiple sentences and entities, while the output consists of all entity pairs in the document and their corresponding relation types.

**JE&RE (Entity-Relation Joint Extraction)**  The previously mentioned relation extraction approaches follow a pipeline where entity recognition is performed first, followed by relation classification based on the identified entities. In contrast, JE&RE task differs by requiring the model to extract both entities and relations simultaneously, without dividing the process into two separate tasks. In this work, we conduct experiments using the FewRel and NYT datasets.

## F.2  EKG

**SED (Sentence-level Event Detection)**  Event detection (ED) aims to identify the events mentioned in a given text and recognize their characteristics, such as event type, participants, time, and other relevant attributes. SED is a specific form of ED, where the task requires the model to detect events within individual sentences. In this work, we utilize the ACE2005 dataset for training and testing the model.

**DED (Document-level Event Detection)**  DED aims to identify multiple events within a document and extract relevant information, such as participants, triggers, and other attributes. Since these events may be distributed across different sentences, DED requires the model to have cross-sentence contextual understanding, making it more complex and enriched compared to sentence-level tasks. In this work, we use the WIKIEVENTS dataset, leveraging Wikipedia entries as events to train and test the model.

**DEAE(Document-level Event Argument Extraction)**  DEAE is a task designed to extract argumentative material from a full document, requiring the identification of arguments in a relationship and the extraction of the relations between arguments and events. In our work, we train and test the model using the WIKIEVENTS and RAMS datasets, where the RAMS dataset includes a rich set of argument types and deals with the relations of argument elements between different sentences.

**ETRE (Event Temporal Relation Extraction)**  ETRE aims to extract events mentioned in a text and determine the temporal order in which these events occur. In our experiments, we use the MATRES, ESL, TB-Dense, Causal-TB, MAVEN-ERE, and TCR datasets for training and testing the model. Notably, the TCR dataset, as an **OOD dataset**, is only used for testing and not for training.

**ECRE (Event Causal Relation Extraction)**  ECRE aims to identify and extract causal relationships between different events in a text. In our work, we use the ESL and MAVEN-ERE datasets for training and testing the model. The ESL dataset is further annotated with various types of causal relationships between events, including direct causality, indirect causality, and opposition relationships. Additionally, during testing, we employ the Causal-TB dataset as an **OOD dataset**, which is only used for testing and not for training.

**ESRE (Event Subevent Relation Extraction)**  In complex texts, events often do not exist independently but can exhibit hierarchical structures, where one event may be the cause, effect, or sub-event of another. ESRE aims to identify these hierarchical relationships between events to achieve a more comprehensive understanding of the event timeline and causal chains. The input to this task is typically a text containing multiple events, and the output is pairs of events along with their hierarchical relationship labels, such as parent event and child event, causal relation, and parallel relation. In this work, we use the HiEve and MAVEN-ERE datasets for model training and testing.

## F.3  CKG

**NER (Named Entity Recognition)**  NER aims to identify entities with specific semantic meanings from a text and classify them into predefined categories, such as person names, locations, organizations, dates, times, and numerical values. Given a natural language text as input, the output consists of the extracted named entities and their corresponding categories. NER plays a critical role in the construction of knowledge graphs by recognizing entities in the text and linking them to existing entity nodes in the knowledge graph, facilitating the automated development and expansion of the graph.

In this work, we use the CoNLL dataset for training and testing the NER task.

**AG (Abstract Generation)**   AG aims to compress a lengthy input text into a concise and accurate abstract while retaining key information and themes. Since CKG can provide rich background and relational information, we employ a CKG-based abstraction task. For this purpose, we train and test the model using the CNNDM and XSum datasets, with the **ROUGE-L** percentage metric used as the evaluation criterion.

**LI (Language Inference)**   The task of LI aims to establish an understanding of relationships between sentences. The core objective of this task is to determine whether a given pair of sentences exhibits entailment, contradiction, or neutrality. Typically, the input consists of a pair of texts, and the output indicates whether the relationship between the two sentences is entailment, contradiction, or neutral. In this work, we use two specialized datasets in the field of natural language inference, the SNLI and MNLI datasets, for training and testing the model.

**TC (Text Classification)**   TC task aims to automatically assign textual data to one or more predefined categories. Given a text as input, the output is typically the predicted category or categories corresponding to the input text. In this work, we use the R8 and R52 datasets for model training and testing, with R8 serving as an **OOD dataset** that is used only for training and not for testing.

## F.4   Counter

**NLG (Natural Language Generation)**   NLG aims to generate natural language text in a predefined format or structure based on specific input information or structure. Unlike traditional free-text generation, the structured text generation task emphasizes the structure and accuracy of the information in the output. The input can take various forms of structured data, such as knowledge graphs, tables, or tuples, and the output is typically a coherent piece of text that adheres to the predetermined structure. In this work, we use the WebNLG dataset, a typical dataset in this domain, for model training and testing. Specifically, we employ the **ROUGE-L** percentage metric as the evaluation criterion.