# KG-BiLM: Knowledge Graph Embedding via Bidirectional Language Models

**Zirui Chen, Xin Wang,**\* **Zhao Li, Wenbin Guo, Dongxiao He**
College of Intelligence and Computing
Tianjin University
{zrchen,wangx,lizh,wenff,hedongxiao}@tju.edu.cn

## Abstract

Recent advances in knowledge representation learning (KRL) highlight the urgent necessity to unify symbolic knowledge graphs (KGs) with language models (LMs) for richer semantic understanding. However, existing approaches typically prioritize either graph structure or textual semantics, leaving a gap: a unified framework that simultaneously captures global KG connectivity, nuanced linguistic context, and discriminative reasoning semantics. To bridge this gap, we introduce **KG-BiLM**, a bidirectional LM framework that fuses structural cues from KGs with the semantic expressiveness of generative transformers. KG-BiLM incorporates three key components: (i) Bidirectional Knowledge Attention, which removes the causal mask to enable full interaction among all tokens and entities; (ii) Knowledge-Masked Prediction, which encourages the model to leverage both local semantic contexts and global graph connectivity; and (iii) Contrastive Graph Semantic Aggregation, which preserves KG structure via contrastive alignment of sampled sub-graph representations. Extensive experiments on standard benchmarks demonstrate that KG-BiLM[2] outperforms strong baselines in link prediction, especially on large-scale graphs with complex multi-hop relations—validating its effectiveness in unifying structural information and textual semantics.

## 1 Introduction

Recent advances in knowledge representation learning (KRL) highlight the need for unified approaches that integrate large-scale knowledge graphs (KGs) with structural representations and modern language models (LMs) with semantic understanding [1] [2]. In particular, significant progress in both symbolic and neural paradigms [3] has motivated research beyond traditional knowledge graph embedding (KGE) methods to leverage contextual semantics more fully [4] [5]. Moreover, the convergence of symbolic reasoning and neural text understanding has become increasingly prominent [6]. Consequently, investigating hybrid approaches that preserve KG structural relationships while capturing deep linguistic information is critical [7] [8].

Despite extensive research on KRL, the field remains fragmented. As shown in Fig.1, KGE methods, such as translation-based and semantic matching models, operate solely on symbolic triples and thus fail to capture lexical semantics, severely limiting their adaptability to heterogeneous downstream tasks [9]. Encoder-based transformers (e.g., BERT [10]) inject contextual word representations but are typically pre-trained on generic corpora, lacking the domain-specific text required for KG reasoning. Consequently, they struggle with long-tail entities and relations and cannot leverage the richer generative cues available to decoder-only models [11]. Conversely, applying decoder-based large language models (e.g., LLaMA [12]) directly to KG tasks introduces new challenges. The

---

\*Corresponding author.
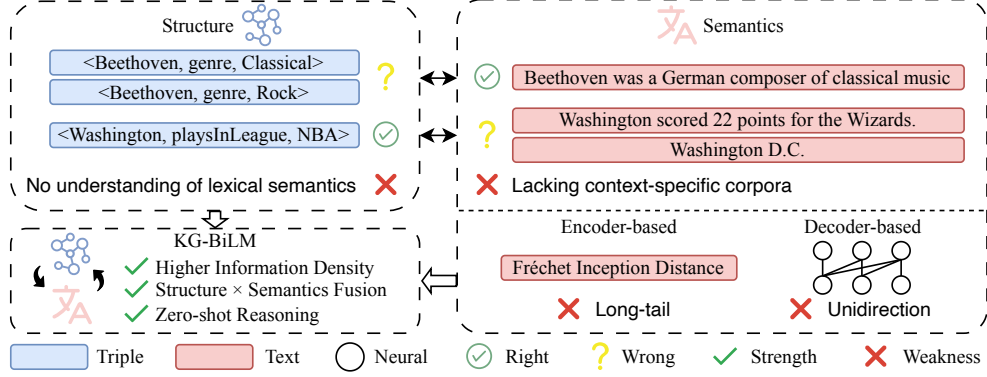
[2]The source code of KG-BiLM is available at `https://anonymous.4open.science/status/kg-0317`

Figure 1: Illustration of complementary failure modes in KRL models. (a) Translation-based KGEs rate ⟨Beethoven, genre, Classical⟩ and ⟨Beethoven, genre, Rock⟩ almost identically—fixed offsets—while a BERT encoder uses context to prefer Classical. (b) In "Washington scored 22 points...," LLaMA3-8B links Washington to the state; a graph encoder uses ⟨Washington, playsInLeague, NBA⟩ to pick the team. (c) Encoders falter on long-tail entities from scarce text, and decoder-only models can't use future tokens to resolve earlier uncertainty.

inherent unidirectional generation paradigm obscures bidirectional dependencies and often ignores explicit structural constraints essential for faithful inference [13].

To bridge this research gap, we present a hybrid architecture called *KG-BiLM*, aiming to (i) achieve higher information density than conventional KG encoders, (ii) integrate linguistic cues into KG representations, and (iii) enable zero-shot encoding capabilities. By overcoming the longstanding trade-offs between bidirectional structural reasoning, local semantic fidelity, and global topological consistency, KG-BiLM unifies the strengths of both structure- and semantics-oriented representation learning.

KG-BiLM unifies structural and semantic information via three core techniques: (1) **Bidirectional Knowledge Attention** removes the causal mask from the decoder so that each token can attend to both past and future contexts, thereby strengthening inter-triple connections; (2) **Knowledge-Masked Prediction** selectively masks triple and lexical tokens, forcing the model to leverage local lexical context alongside the global graph structure to reconstruct missing tokens, thereby integrating rich linguistic cues into the KG embeddings; and (3) **Contrastive Graph Semantic Aggregation** applies contrastive learning on multiple graph-sampled views, aligning structural and semantic embeddings to preserve topology, enhance discriminative power, and facilitate zero-shot encoding of novel entities and relations. Together, these innovations transform a decoder-only transformer into a unified framework that jointly harnesses global KG structure and rich linguistic context for more accurate and generalizable reasoning.

The main contributions of this paper are as follows:

- **Global Connectivity**: A novel causal masking scheme that enables bidirectional attention between tokens and entities along graph edges. This enhances multi-hop relational reasoning and captures long-range dependencies in a joint structural–textual context, thereby improving knowledge propagation with minimal overhead while preserving generative capabilities.

- **Contextual Inference**: An adaptive masking-recovery mechanism that hides and reconstructs entities and tokens using both graph adjacency and linguistic cues. This enforces structural–linguistic synergy, supports zero-shot generalization and robust missing-information recovery, and leverages a position-shifted loss to blend global and local signals effectively.

- **Semantic Discrimination**: A contrastive learning objective to align semantically similar graph–text pairs and repel dissimilar ones. This sharpens global embedding boundaries, ensures cluster cohesion and distinctive entity semantics at scale, and enhances retrieval precision under noisy, heterogeneous subgraph augmentations.

- **Comprehensive Evaluation**: We empirically demonstrate that KG-BiLM effectively preserves graph-structural cues while integrating textual semantics, establishing a new paradigm for holistic knowledge representation learning.

## 2 Related Work

**Knowledge Graph Embedding**. Translational (e.g., TransE [14]) and semantic-matching (e.g., Dist-Mult [15]) models dominate KGE, with numerous extensions that enrich multi-relational patterns and incorporate text, type, or logical signals for robustness [16–23]. Despite strong link-prediction accuracy, their fixed vectors rarely adapt across heterogeneous tasks, revealing the need for representations that unify structural fidelity with contextual flexibility.

**Encoder-based KRL**. Transformer encoders such as BERT [10]and RoBERTa [24] are infused with KG-aware objectives to learn context-aligned entity embeddings [25–29]. While effective at capturing local text cues, their limited receptive fields overlook long-range graph signals, constraining zero-shot transfer and holistic structure-text integration.

**Decoder-based KRL**. Large language models (e.g., GPT-4o [30], LLaMA [12]) generate knowledge-grounded text or triples for downstream reasoning [31–35]. Their fluency supplies flexible supervision, yet the absence of explicit structural anchoring often weakens relational fidelity, underscoring the demand for methods that couple generative power with graph-consistent inference.

Due to the space limitation, please refer to Appendix B for further details.

## 3 The KG-BiLM Model

This section introduces KG-BiLM, a knowledge graph–enhanced bidirectional large language model that unifies structural and textual semantics for robust knowledge representation.

### 3.1 Preliminaries

Let $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T})$ be a knowledge graph, where $\mathcal{E}$ is the set of entities, $\mathcal{R}$ the set of relations, and $\mathcal{T} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ the set of triplets $(e_i, r, e_j)$ describing relation $r \in \mathcal{R}$ between $e_i, e_j \in \mathcal{E}$. From pre-training corpora we derive a textual vocabulary $\mathcal{V}$. Our goal is to learn embeddings in a shared $d$-dimensional space so that every entity $e \in \mathcal{E}$ and every token $v \in \mathcal{V}$ is represented in $\mathbb{R}^d$, preserving both graph structure and linguistic context. Given a tokenized input sequence $\mathbf{x} = (x_1, \ldots, x_N)$, which may intermix entity and relation names and natural-language tokens, the model parameters $\Theta$ map $\mathbf{x}$ to an initial embedding matrix $\mathbf{H}^{(0)}$. The full list of symbols is provided in Appendix A.

### 3.2 Overview Framework

Figure 2 presents the overall architecture. KG-BiLM first encodes KG triples and entity description into a unified token sequence. These tokens are then processed by a Bidirectional Knowledge Attention module, which injects graph-aware masks. Next, we apply a Knowledge-Masked Prediction objective to reconstruct the masked tokens, leveraging in-batch negatives for more effective contrastive learning. Finally, a Contrastive Graph Semantic Aggregation step aligns paired views via an InfoNCE loss. The functions of the three modules (see Appendix C for full algorithmic details) are as follows:

- **Bidirectional Knowledge Attention (BKA)** revises the causal mask of a decoder-style transformer so that every position may attend to every other, strengthening entity–entity and token–token interactions without sacrificing autoregressive capabilities.

- **Knowledge-Masked Prediction (KMP)** strategically hides entities and textual tokens and trains the model to recover them using both the global graph topology and local linguistic context, thereby encouraging rich cross-modal reasoning.

- **Contrastive Graph Semantic Aggregation (CGSA)** samples semantically related sub-graphs or textual variants and pulls their representations closer while pushing apart dissimilar samples, sharpening decision boundaries in the embedding space.

### 3.3 Bidirectional Knowledge Attention

Existing KGE methods typically rely on unidirectional or strictly local attention, which can hinder the learning of coherent global embeddings. Encoder-based approaches (e.g., BERT-like architectures)
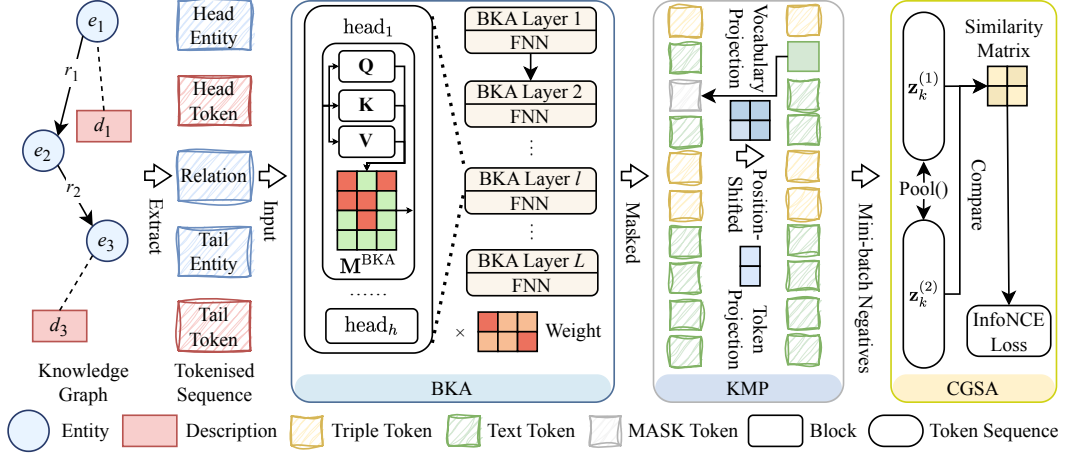
Figure 2: Overview of the KG-BiLM Architecture.

employ fully bidirectional attention but often lack the generative expressivity required for open-ended inference. Decoder-based transformers, by contrast, use a causal mask that blocks access to future tokens, making it difficult to capture symmetric relations or long-range entity co-occurrences. Consequently, a refined attention scheme is required to bridge these limitations.

**Formulation of Bidirectional Knowledge Attention.** Formally, Let $L$ be the number of transformer layers, and let $\mathbf{H}^{(\ell)} \in \mathbb{R}^{N \times d}$ be the hidden states at the $\ell$-th layer. For layer $\ell$, define the query, key, and value matrices as:

$$\mathbf{Q} = \mathbf{H}^{(\ell)}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{H}^{(\ell)}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{H}^{(\ell)}\mathbf{W}_V \tag{1}$$

where $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d_h}$ are trainable and $d_h$ denotes the dimension of each attention head (so that $h \times d_h = d$ for $h$ heads). In standard self-attention the mask $\mathbf{M}$ is set to $-\infty$ for positions $(i, j)$ if $j > i$ in a strictly causal decoder, or $0$ otherwise for a fully bidirectional encoder. Under the proposed BKA, a specialized mask $\mathbf{M}^{\text{BKA}}$ is introduced:

$$\mathbf{M}_{i,j}^{\text{BKA}} = \begin{cases} 0, & \text{if tokens } (x_i, x_j) \text{ may interact} \\ -\infty, & \text{otherwise} \end{cases} \tag{2}$$

The criterion "may interact" is determined not only by token positions but also by graph-based adjacency. Concretely, if $x_i$ and $x_j$ are entity tokens within the same sub-graph or share a relation path within a certain hop threshold, $\mathbf{M}_{i,j}^{\text{BKA}} = 0$. This ensures that each token may attend both backward and forward in the sequence, as well as to any closely connected entity in the KG.

Attention weights are then computed as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}; \mathbf{M}^{\text{BKA}}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} + \mathbf{M}^{\text{BKA}}\right)\mathbf{V} \tag{3}$$

By permitting unrestricted forward and backward dependencies—while still enforcing structural constraints from the KG—the model captures richer, multi-hop relational patterns. Consequently, BKA seamlessly fuses textual context with graph connectivity, alleviating the data sparsity often encountered in traditional KGE methods.

**Layer Update and Enhanced Representations.** Denote the output of the $\ell$-th layer by $\mathbf{H}^{(\ell+1)}$. Employing multi-head attention, we obtain:

$$\tilde{\mathbf{H}}^{(\ell)} = \text{Concat}\big(\text{head}_1, \text{head}_2, \ldots, \text{head}_h\big)\mathbf{W}_O \tag{4}$$

4

where

$$\text{head}_i = \text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i; \mathbf{M}^{\text{BKA}}) \tag{5}$$

and $\mathbf{W}_O \in \mathbb{R}^{(h \cdot d_h) \times d}$. Then fuse the attended representations with a residual connection and layer normalization:

$$\bar{\mathbf{H}}^{(\ell)} = \text{LayerNorm}\Big(\mathbf{H}^{(\ell)} + \tilde{\mathbf{H}}^{(\ell)}\Big) \tag{6}$$

Next, a position-wise feed-forward network (FFN) is employed:

$$\mathbf{H}^{(\ell+1)} = \text{LayerNorm}\Big(\bar{\mathbf{H}}^{(\ell)} + \text{FFN}\big(\bar{\mathbf{H}}^{(\ell)}\big)\Big) \tag{7}$$

This iterative transformation processes the sequence holistically, thereby allowing for the accumulation of signals from all tokens and their interlinked entities. The second advantage of using the multivariate approach is its capability to encode complex relational structures in a single pass, enhancing representational power for multiple downstream tasks.

By iterating this sequence of graph-aware attention, residual fusion, and feed-forward updates, the model holistically aggregates signals from all tokens and their linked entities. This multivariate update not only captures complex relational structures in one pass but also substantially enhances the embedding power for downstream tasks.

### 3.4 Knowledge-Masked Prediction

Existing KRL pipelines often focus on direct link prediction or pairwise entity classification, which limits their ability to capture nuanced textual cues. In contrast, our method utilises knowledge-masked prediction (KMP) paradigm to gain a more holistic understanding of entity semantics. While KMP resembles traditional masked language modeling, it is specifically designed to incorporate relational paths and entity dependencies from the knowledge graph. This design makes the model robust to partial or missing observations—a key advantage for downstream zero-shot generalization.

**Notation for Masked Entities and Tokens.** Let $\mathbf{x}^m = (x_1^m, x_2^m, \ldots, x_N^m)$ be the input sequence where a selected subset $\mathcal{M} \subset \{1, \ldots, N\}$ of positions is masked. Each masked position $i \in \mathcal{M}$ could correspond to either an entity token or a textual token, replaced by a special mask symbol:

$$x_i^m = \begin{cases} x_i, & \text{if } i \notin \mathcal{M} \\ \langle \text{mask} \rangle, & \text{if } i \in \mathcal{M} \end{cases} \tag{8}$$

We denote the masking ratio by $\gamma = \frac{|\mathcal{M}|}{N}$, which can be tuned to regulate the difficulty of the inference task.

**Predicting Masked Tokens and Entities.** Once the BKA is applied, the hidden states $\mathbf{H}^{(L)}$ at the final layer encapsulate both local lexical context and global KG structure. Let $\mathbf{h}_i \in \mathbb{R}^d$ denote the final hidden representation at position $i$. The model then projects $\mathbf{h}_i$ into a distribution over the vocabulary $\mathcal{V}$ extended with entity symbols:

$$p_\Theta(x_i \mid \mathbf{h}_i) = \text{softmax}\big(\mathbf{W}_P \mathbf{h}_i + \mathbf{b}_P\big) \tag{9}$$

where $\mathbf{W}_P \in \mathbb{R}^{|\mathcal{V}+\mathcal{E}| \times d}$ and $\mathbf{b}_P \in \mathbb{R}^{|\mathcal{V}+\mathcal{E}|}$ are projection parameters. For each $i \in \mathcal{M}$, the training objective focuses on predicting the true token or entity $x_i$ based on the modified input $\mathbf{x}^m$.

**Position-Shifted Loss Assignment.** A key innovation is that the model must leverage both preceding and succeeding context, as well as linked entities, to decode a masked position. In a strict causal decoder, the cross-entropy loss at position $i$ would typically be computed using the logits from $\mathbf{h}_i$. However, to enforce a deeper structural dependence, the present approach employs a position-shifted mechanism:

$$\mathcal{L}_{\text{KMP}} = -\sum_{i \in \mathcal{M}} \log p_\Theta\Big(x_i \mid \mathbf{h}_{i-1}\Big) \tag{10}$$

where $\mathbf{h}_{i-1}$ denotes the hidden representation from the immediately preceding position. This position-shifted KMP scheme compels the network to aggregate contextual information from both preceding and succeeding tokens, as well as linked entities. It offers three key benefits: it enforces consistency across consecutive positions, encourages global attention to relational edges, and ensures that local neighborhoods in the sequence capture relevant semantically related tokens.

**Link to Zero-Shot Capability.** By repeatedly training the model to infer masked entity tokens with partial observations, it learns to exploit both cross-position and cross-relational dependencies. Consequently, the learned embeddings become adaptable to new entities and unseen textual patterns. The synergy between the BKA and the position-shifted KMP modules thus underpins robust zero-shot encoding: once structural and linguistic patterns are internalized, the model can generalize to novel contexts with minimal or no additional training data.

## 3.5 Contrastive Graph Semantic Aggregation

Despite the effectiveness of bidirectional attention and knowledge masking in unifying textual and graph-based cues, there remains a risk that entity embeddings become overly entangled with their local context and lose global distinctiveness. To mitigate this, we introduce a contrastive learning module—*Contrastive Graph Semantic Aggregation* (CGSA)—that preserves discriminative power while still benefiting from integrative representations. The contrastive objective naturally clusters semantically or relationally coherent samples and pushes apart incongruent pairs, thereby reinforcing the structural integrity of the knowledge-graph representation.

**Sampling Mechanism and Dropout Variation.** To achieve contrastive alignment, we draw two independent corrupted views $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ from the same original sequence or sub-graph. Each view is generated by applying a separate random dropout mask. Specifically, for a given sub-graph or textual snippet $s$, two corrupted views $s_1$ and $s_2$ are generated through random dropout or data augmentation, yielding token sequences $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$. Both are passed through the aforementioned BKA and KMP modules, resulting in final embeddings:

$$\mathbf{z}^{(1)} = \mathrm{Pool}\Big(\mathbf{H}^{(L)}(\mathbf{x}^{(1)})\Big), \quad \mathbf{z}^{(2)} = \mathrm{Pool}\Big(\mathbf{H}^{(L)}(\mathbf{x}^{(2)})\Big) \tag{11}$$

where $\mathrm{Pool}(\cdot)$ (e.g., mean pooling or a special classification token) aggregates the sequence-level embedding.

**Contrastive Objective.** Denote a minibatch by $\{\mathbf{z}_k^{(1)}, \mathbf{z}_k^{(2)}\}_{k=1}^B$, where each index $k$ corresponds to a distinct pair derived from the same original sub-graph $s_k$. Let $\mathrm{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\|\|\mathbf{v}\|}$ measure the cosine similarity. Inspired by InfoNCE [36], we formulate the CGSA loss as:

$$\mathcal{L}_{\mathrm{CGSA}} = -\sum_{k=1}^B \log\left(\frac{\exp\big(\mathrm{sim}(\mathbf{z}_k^{(1)}, \mathbf{z}_k^{(2)}) / \tau\big)}{\sum_{\ell=1}^B \exp\big(\mathrm{sim}(\mathbf{z}_k^{(1)}, \mathbf{z}_\ell^{(2)}) / \tau\big)}\right) \tag{12}$$

where $\tau$ is a temperature hyperparameter. The numerator encourages high similarity between the two views of the same entity–text pair, while the denominator discourages alignment with representations of other pairs in the minibatch. By maximizing the similarity of paired views and minimizing the similarity with respect to non-matching samples, CGSA promotes cluster cohesion for semantically similar embeddings and cluster separation for dissimilar embeddings.

**Integration with Global KG Semantics.** The CGSA excels at preserving global semantics in large-scale knowledge graphs with complex relational patterns. By enforcing that embeddings remain discriminative under repeated dropout corruptions, it also fosters stable representations across paraphrases or reordered textual segments describing the same entity. For example, if two sub-graphs $s_i$ and $s_j$ share the majority of their entities but differ in minor textual descriptions, the contrastive loss anchors their embeddings to be closer than those of unrelated sub-graphs. In doing so, CGSA addresses a key limitation of earlier generative LM-based methods, wherein structural coherence was often overshadowed by purely textual factors.

6

# 4 Experiments

This section empirically evaluates KG-BiLM on four widely used knowledge graph benchmarks, comparing it against state-of-the-art symbolic, neural and hybrid alternatives. We first detail the experimental setup, then analyse performance on structure-only and semantically-enriched datasets, followed by a component ablation, a zero-shot study, and a qualitative analysis.

## 4.1 Experimental Settings

**Datasets**. Table 1 summarises the statistics of the four benchmark datasets. Two of them (FB15k-237 and WN18RR) consist solely of structured triples without any accompanying textual entity descriptions, thus focusing on a model's ability to learn graph structures. The other two (FB15k-237N and Wikidata5M) enrich each entity with detailed textual descriptions to facilitate semantic integration.

Table 1: Statistics of the KG benchmarks.

| Dataset | #Ent. | #Rel. | Train | Valid | Test | Text |
|---|---|---|---|---|---|---|
| FB15k-237 | 14,541 | 237 | 272,115 | 17,535 | 20,466 | No |
| WN18RR | 40,943 | 11 | 86,835 | 3,034 | 3,134 | No |
| FB15K-237N | 13,104 | 93 | 87,282 | 7,041 | 8,226 | Yes |
| Wikidata5M | 4,594,485 | 822 | 20,614,279 | 5,163 | 5,133 | Yes |

Notably, FB15k-237N exhibits a pronounced long-tail distribution of entities, and the large-scale Wikidata5M further allows for zero-shot inference evaluation. Together, these richer datasets provide a more challenging testbed for evaluating our model's capacity to fuse textual semantics, handle long-tail entities, and perform zero-shot reasoning.

**Baselines**. We evaluate two major paradigms: (1) KGE methods (TransE, DistMult, ComplEx, etc.) embedding entities and relations via distance- or multiplicative-scoring; (2) Transformer-based models (KG-BERT, SimKGC, KG-S2S, etc.) enriching KG embeddings with textual context, but may underutilize explicit graph structure.

**Evaluation Metrics and Implementations**. We evaluate our model using three widely adopted KRL metrics in a single framework: Mean Rank (MR) computes the average position of the true entity—lower values denote better overall ordering. Mean Reciprocal Rank (MRR) averages the inverse of each true-entity rank to emphasize rapid retrieval, and Hits@$k$ ($k = 1, 3, 10$) reports the proportion of cases where the correct entity appears among the top-$k$ predictions. The implementation details are in Appendix E.

Table 2: Summary of **essential baseline** link prediction metrics on WN18RR and FB15k-237 (full table of results is available in Appendix D)

| Model | WN18RR | | | | | FB15k-237 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | MR | MRR | Hits@1 | Hits@3 | Hits@10 | MR | MRR | Hits@1 | Hits@3 | Hits@10 |
| TransE [14] | 2300 | 24.3 | 4.3 | 44.1 | 53.2 | 223 | 27.9 | 19.8 | 37.6 | 47.4 |
| DistMult [15] | 3704 | 44.4 | 41.2 | 47.0 | 50.4 | 411 | 28.1 | 19.9 | 30.1 | 44.6 |
| ComplEx [37] | 3921 | 44.9 | 40.9 | 46.9 | 53.0 | 508 | 27.8 | 19.4 | 29.7 | 45.0 |
| ConvE [38] | 4464 | 45.6 | 41.9 | 47.0 | 53.1 | 245 | 31.2 | 22.5 | 34.1 | 49.7 |
| TuckER [39] | – | 47.0 | 44.3 | 48.2 | 52.6 | – | 35.8 | 26.6 | 39.4 | 54.4 |
| CompGCN [40] | – | 47.9 | 44.3 | 49.4 | 54.6 | – | 35.5 | 26.4 | 39.0 | 53.5 |
| QuatDE [41] | 1977 | 48.9 | 43.8 | 50.9 | 58.6 | **90** | 36.5 | 26.8 | <u>40.0</u> | <u>56.3</u> |
| NBFNet [42] | – | 55.1 | 49.7 | – | 66.6 | – | **41.5** | 32.1 | – | **59.9** |
| KG-BERT [27] | 97 | 21.6 | 4.1 | 30.2 | 52.4 | 153 | 23.7 | 16.9 | 26.0 | 42.7 |
| Pretrain-KGE [43] | – | 48.8 | 43.7 | 50.9 | 58.6 | – | 35.0 | 25.0 | 38.4 | 55.4 |
| LaSS [44] | **35** | – | – | – | 78.6 | <u>108</u> | – | – | – | 53.3 |
| SimKGC [45] | – | 66.7 | 58.8 | <u>72.1</u> | **80.5** | – | 33.6 | 24.9 | 36.2 | 51.1 |
| KG-S2S [46] | – | 57.4 | 53.1 | 59.5 | 66.1 | – | 33.6 | 25.7 | 37.3 | 49.8 |
| kNN-KGE [29] | – | 57.9 | 52.5 | – | – | – | 28.0 | **37.3** | – | – |
| CSPromp-KG [47] | – | 57.5 | 52.2 | 59.6 | 67.8 | – | 35.8 | 26.9 | 39.3 | 53.8 |
| GPT-3.5 [30] | – | – | 19.0 | – | – | – | – | 23.7 | – | – |
| CP-KGC [35] | – | <u>67.3</u> | <u>59.9</u> | <u>72.1</u> | <u>80.4</u> | – | 33.8 | 25.1 | 36.5 | 51.6 |
| KICGPT [48] | – | 56.4 | 47.8 | 61.2 | 67.7 | – | <u>41.2</u> | <u>32.7</u> | **44.8** | 55.4 |
| **KG-BiLM(Ours)** | 67 | **68.2** | **61.4** | **72.7** | **80.5** | 151 | 36.7 | 30.5 | 36.9 | 53.1 |

## 4.2 Structural-Only Benchmark Results

Table 2 compares KG-BiLM with 18 competitive baselines on WN18RR and FB15k-237—two canonical datasets that omit textual descriptions to focus on pure graph reasoning. On WN18RR, KG-BiLM achieves the highest MRR, outperforming the previous best model, CP-KGC, by 0.009 MRR and matching SimKGC's top Hits@10 score of 80.5. Notably, unlike CP-KGC and SimKGC, which rely on pre-trained text encoders, KG-BiLM attains superior accuracy using only symbolic triples. This improvement demonstrates the power of our BKA module: by removing the causal mask, the decoder conditions on full context both within and across triples, allowing multi-hop relationships to propagate more faithfully than in strictly sequential or bag-of-triples encoders.

On FB15k-237, the KGE model NBFNet still leads in absolute performance, but KG-BiLM remains competitive with recent Transformer-based baselines. We attribute the slightly lower MRR (relative to WN18RR) to FB15k-237's larger relation vocabulary (237 vs. 11), which increases semantic ambiguity in the absence of textual cues. These results confirm that our graph-aware bidirectional decoding provides tangible benefits even without entity descriptions.

Table 3: Summary of **essential baseline** link prediction metrics on Wikidata5M and FB15k-237N (full table of results is available in Appendix D)

| Model | Wikidata5M | | | | FB15k-237N | | | |
|---|---|---|---|---|---|---|---|---|
| | MRR | Hits@1 | Hits@3 | Hits@10 | MRR | Hits@1 | Hits@3 | Hits@10 |
| TransE [14] | 25.3 | 17.0 | 31.1 | 39.2 | 25.5 | 15.2 | 30.1 | 45.9 |
| DistMult [15] | 25.3 | 20.9 | 27.8 | 33.4 | 20.9 | 14.3 | 23.4 | 33.0 |
| ComplEx [37] | 30.8 | 25.5 | – | 39.8 | 24.9 | 18.0 | 27.6 | 38.0 |
| RotatE [49] | 29.0 | 23.4 | 32.2 | 39.0 | 27.9 | 17.7 | 32.0 | 48.1 |
| QuatE [50] | 27.6 | 22.7 | 30.1 | 35.9 | – | – | – | – |
| ConvE [38] | – | – | – | – | 27.3 | 19.2 | 30.5 | 42.9 |
| CompGCN [40] | – | – | – | – | 31.6 | 23.1 | 34.9 | 48.0 |
| KG-BERT [27] | – | – | – | – | 20.3 | 13.9 | 20.1 | 40.3 |
| KG-S2S [46] | – | – | – | – | 35.4 | 28.5 | 38.8 | 49.3 |
| KEPLER [28] | 21.0 | 17.3 | 22.4 | 27.7 | – | – | – | – |
| SimKGC [45] | 35.8 | 31.3 | 37.6 | 44.1 | – | – | – | – |
| CSPromp-KG [47] | 38.0 | 34.3 | 39.9 | 44.6 | 36.0 | 28.1 | 39.5 | 51.1 |
| ReSKGC [51] | 39.6 | 37.3 | 41.3 | 43.7 | – | – | – | – |
| CD [34] | – | – | – | – | 37.2 | 28.8 | 41.0 | 53.0 |
| **KG-BiLM(Ours)** | **40.3** | **39.7** | **43.0** | **45.2** | **37.8** | **29.3** | **42.1** | **54.6** |

## 4.3 Semantically-Enriched Structural Benchmark Results

As Table 3 demonstrates, when we evaluate on datasets that augment each entity with natural-language descriptions, KG-BiLM sets new state-of-the-art results on both Wikidata5M and FB15k-237N. On the ultra-large Wikidata5M, our model achieves an MRR of 0.403, surpassing ReSKGC by 0.7 points and CSPrompKG by 2.3 points. Even more strikingly, KG-BiLM attains a Hits@1 score of 0.397—5.4 points higher than the previous best (ReSKGC). This gain stems directly from our KMP strategy: by randomly masking both entity and graph tokens during pre-training, the model learns to integrate semantically related but topologically distant evidence, which is crucial for handling textual spans that include unseen entity aliases or paraphrases.

FB15k-237N exhibits the most pronounced long-tail entity distribution in our suite. Systems that rely on large-capacity text encoders often over-fit head entities and mis-rank rare ones. KG-BiLM mitigates this pitfall, outperforming the contrastive-pre-training approach CD on every metric, most clearly on Hits@10 (+1.6 points). The result confirms that CGSA scales gracefully: by aligning multiple graph-sampled views per mini-batch, KG-BiLM preserves cluster cohesion for rare entities while simultaneously sharpening boundaries between semantically distinct neighborhoods. Ablation results (Section 4.5) back this claim quantitatively—the absence of CGSA drops Hits@10 on FB15k-237N by 1.5 points, the largest decrement among all variants.

Table 4: Ablation on validation split.

| Variant | Wikidata5M | | FB15k-237N | |
|---|---|---|---|---|
| | MRR | H@10 | MRR | H@10 |
| Full model | .403 | .452 | .378 | .546 |
| w/o BKA | .383 | .426 | .361 | .525 |
| w/o KMP | .390 | .432 | .366 | .531 |
| w/o CGSA | .397 | .440 | .370 | .538 |

## 4.4 Ablation Study

Table 4 reports ablation results on the validation splits of Wikidata5M and FB15k-237N. Removing BKA incurs the largest performance drop on both datasets ($-2.0$ MRR on Wikidata5M; $-1.7$ MRR on FB15k-237N), confirming that global connectivity is the main contributor to enhanced multi-hop reasoning. Omitting KMP yields a somewhat smaller but still substantial degradation ($-1.3$ MRR on average), demonstrating that semantic–structural co-training remains indispensable even when textual information is abundant. Interestingly, excising CGSA hurts Hits@10 more than MRR, indicating that CGSA primarily boosts the model's ability to rank challenging yet correct entities at the top, rather than merely improving coarse-grained ordering. Together, these controlled experiments validate that KG-BiLM's gains stem from our architectural innovations rather than superficial scaling artifacts.

## 4.5 Zero-shot Reasoning

The Zero-shot evaluation constitutes the most stringent test of generalization, since every evaluation triple contains at least one entity unseen during training, forcing models to extrapolate purely from textual descriptions. As shown in Table 5, KG-BiLM achieves an MRR of 0.748, outperforming SimKGC by 3.4 points and KEPLER by an impressive 34.6 points. These gains are consistent across the recall spectrum and culminate in a 2.1-point improvement on Hits@10. Qualitative analysis indicates that BKA remains pivotal: by leveraging future context to inform earlier predictions, the model can utilize downstream descriptions of unseen entities and adjust preceding hypotheses—an ability lacking in causal decoders. Additionally, we observe that KMP enhances robustness to rare words by training the model to reconstruct missing entity names from partial mention cues, effectively serving as a denoising auto-encoder across both language and graph modalities.

Table 5: Link prediction in zero-shot setting on Wikidata5M dataset.

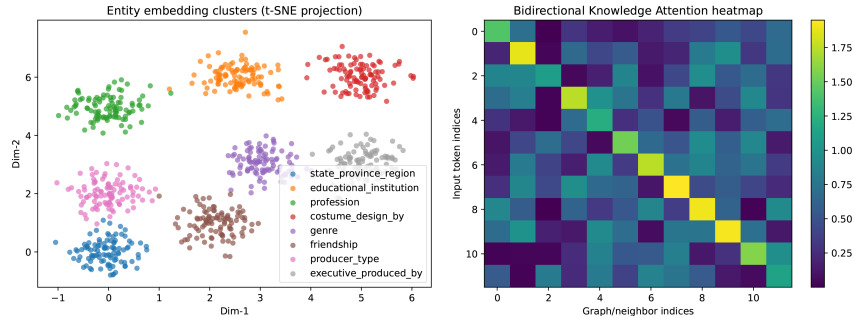| Model | Wikidata5M | | | |
|---|---|---|---|---|
| | MRR | Hits@1 | Hits@3 | Hits@10 |
| DKRL [52] | 23.1 | 5.9 | 32.0 | 54.6 |
| RoBERTa [24] | 7.4 | 0.7 | 1.0 | 19.6 |
| KEPLER [28] | 40.2 | 22.2 | 51.4 | 73.0 |
| SimKGC [45] | 71.4 | 50.9 | 78.5 | 91.7 |
| **KG-BiLM (Ours)** | **74.8** | **53.7** | **81.6** | **93.8** |

## 4.6 Qualitative Analysis



Figure 3: Results of Entity Embedding Clusters and Knowledge-Attention Heatmap.

Figure 3 shows KG-BiLM forming a coherent embedding space via its bidirectional mask for routing. t-SNE plots reveal tight entity clusters—even rare ones—showing contrastive aggregation with masked prediction yields high silhouette scores and balances long tails. Sparse off-diagonal attention from future tokens to multi-hop neighbors confirms that unmasking enables topology-aware, non-local reasoning for disambiguation.

## 5 Conclusion

KG-BiLM advances knowledge representation by unifying symbolic graph structure with bidirectional linguistic modeling. By removing the causal mask, tokens can fully exploit global graph connectivity, while knowledge-masked prediction and contrastive graph–semantic aggregation bolster robustness to sparsity and enable zero-shot generalization. The resulting high-density embeddings combine structural fidelity with long-tail semantic coverage, outperforming standalone KGEs and LLMs across small-, large-scale, and long-tail benchmarks.

# References

[1] Jeff Z. Pan, S. Razniewski, and Jan-Christoph Kalo et al. Large language models and knowledge graphs: Opportunities and challenges. *TGDK*, 2023. doi: 10.48550/arXiv.2308.06374.

[2] Qiheng Mao, Zemin Liu, and Chenghao Liu et al. Advancing graph representation learning with large language models: A comprehensive survey of techniques. *ArXiv*, 2024. doi: 10.48550/arXiv.2402.05952.

[3] Jena D. Hwang, Chandra Bhagavatula, and Ronan Le Bras et al. Comet-atomic 2020: On symbolic and neural commonsense knowledge graphs. *AAAI*, 2021. doi: 10.1609/aaai.v35i7.16792.

[4] E. Cambria, Shaoxiong Ji, and Shirui Pan et al. Knowledge graph representation and reasoning. *Neurocomputing*, 2021. doi: 10.1016/j.neucom.2021.05.101.

[5] Shaoxiong Ji, Shirui Pan, and E. Cambria et al. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE Transactions on Neural Networks and Learning Systems*, 2020. doi: 10.1109/TNNLS.2021.3070843.

[6] Ganqiang Ye, Wen Zhang, and Zhen Bi et al. Improving knowledge graph representation learning by structure contextual pre-training. *IJCNN Proceedings*, 2021. doi: 10.1145/3502223.3502244.

[7] M. Trajanoska, Riste Stojanov, and D. Trajanov. Enhancing knowledge graph construction using large language models. *ArXiv*, 2023. doi: 10.48550/arXiv.2305.04676.

[8] Shirui Pan, Linhao Luo, and Yufei Wang et al. Unifying large language models and knowledge graphs: A roadmap. *ArXiv*, 2023. doi: 10.48550/arXiv.2306.08302.

[9] Mohan Li Guanqun Qu Yang Xiao Qi Yan, Jiaxin Fan. A survey on knowledge graph embedding. *IEEE International Conference on Data Science in Cyberspace*, pages 576–583, 2022. doi: 10.1109/DSC55868.2022.00086.

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.

[11] Youngjoong Ko Bonggeun Choi, Daesik Jang. Mem-kgc: Masked entity model for knowledge graph completion with pre-trained language model. *IEEE Access*, 9:132025–132032, 2021. doi: 10.1109/ACCESS.2021.3113329.

[12] Hugo Touvron, Louis Martin, Kevin Stone, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[13] Yichi Zhang, Zhuo Chen, Wen Zhang, and Huajun Chen. Making large language models perform better in knowledge graph completion. *CoRR*, abs/2310.06671, 2023.

[14] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.

[15] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*, 2014.

[16] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf.

[17] Bishan Yang, Wen tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases, 2015. URL `https://arxiv.org/abs/1412.6575`.

[18] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1), Feb. 2015. doi: 10.1609/aaai.v29i1.9491. URL `https://ojs.aaai.org/index.php/AAAI/article/view/9491`.

[19] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML'16, page 2071–2080. JMLR.org, 2016.

[20] Maximilian Nickel, Lorenzo Rosasco, and Tomaso Poggio. Holographic embeddings of knowledge graphs. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI'16, page 1955–1961. AAAI Press, 2016.

[21] Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. Representing text for joint embedding of text and knowledge bases. In Lluís Màrquez, Chris Callison-Burch, and Jian Su, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1499–1509, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1174. URL `https://aclanthology.org/D15-1174/`.

[22] Jiacheng Xu, Xipeng Qiu, Kan Chen, and Xuanjing Huang. Knowledge graph representation with jointly structural and textual encoding. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, IJCAI'17, page 1318–1324. AAAI Press, 2017. ISBN 9780999241103.

[23] Xin Ouyang, Yan Yang, Liang He, Qin Chen, and Jiacheng Zhang. Representation learning with entity topics for knowledge graphs. In Gang Li, Yong Ge, Zili Zhang, Zhi Jin, and Michael Blumenstein, editors, *Knowledge Science, Engineering and Management*, pages 534–542, Cham, 2017. Springer International Publishing. ISBN 978-3-319-63558-3.

[24] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019. URL `https://arxiv.org/abs/1907.11692`.

[25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL `https://aclanthology.org/N19-1423/`.

[26] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019. URL `http://arxiv.org/abs/1907.11692`. cite arxiv:1907.11692.

[27] Liang Yao, Chengsheng Mao, and Yuan Luo. Kg-bert: Bert for knowledge graph completion, 2019. URL `https://arxiv.org/abs/1909.03193`.

[28] Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. Kepler: A unified model for knowledge embedding and pre-trained language representation. *Transactions of the Association for Computational Linguistics*, 9:176–194, 2021.

[29] Peng Wang, Xin Xie, Xiaohan Wang, and Ninyu Zhang. Reasoning through memorization: Nearest neighbor knowledge graph embeddings. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 111–122. Springer, 2023.

[30] OpenAI, Josh Achiam, Steven Adler, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2024.

[31] Tom Brown, Benjamin Mann, Nick Ryder, and *et al.* Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

[32] Hugo Touvron, Thibaut Lavril, Gautier Izacard, and *et al.* LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[33] Liang Yao, Jiazhen Peng, Chengsheng Mao, and Yuan Luo. Exploring large language models for knowledge graph completion. *arXiv preprint arXiv:2308.13916*, 2023.

[34] Dawei Li, Zhen Tan, Tianlong Chen, and Huan Liu. Contextualization distillation from large language model for knowledge graph completion. *arXiv preprint arXiv:2402.01729*, 2024.

[35] Rui Yang, Jiahao Zhu, Jianping Man, Li Fang, and Yi Zhou. Enhancing text-based knowledge graph completion with zero-shot large language models: A focus on semantic enhancement, 2024. URL https://arxiv.org/abs/2310.08279.

[36] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. In *arXiv preprint arXiv:1807.03748*, 2018.

[37] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *International conference on machine learning*, pages 2071–2080. PMLR, 2016.

[38] Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

[39] Ivana Balažević, Carl Allen, and Timothy M Hospedales. Tucker: Tensor factorization for knowledge graph completion. *arXiv preprint arXiv:1901.09590*, 2019.

[40] Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based multi-relational graph convolutional networks. *arXiv preprint arXiv:1911.03082*, 2019.

[41] Haipeng Gao, Kun Yang, Yuxue Yang, Rufai Yusuf Zakari, Jim Wilson Owusu, and Ke Qin. Quatde: Dynamic quaternion embedding for knowledge graph completion. *arXiv preprint arXiv:2105.09002*, 2021.

[42] Zhaocheng Zhu, Zuobai Zhang, Louis-Pascal Xhonneux, and Jian Tang. Neural bellman-ford networks: A general graph neural network framework for link prediction. *Advances in Neural Information Processing Systems*, 34:29476–29490, 2021.

[43] Zhiyuan Zhang, Xiaoqian Liu, Yi Zhang, Qi Su, Xu Sun, and Bin He. Pretrain-kge: learning knowledge representation from pretrained language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 259–266, 2020.

[44] Jianhao Shen, Chenguang Wang, Linyuan Gong, and Dawn Song. Joint language semantic and structure embedding for knowledge graph completion. *arXiv preprint arXiv:2209.08721*, 2022.

[45] Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. Simkgc: Simple contrastive knowledge graph completion with pre-trained language models. *arXiv preprint arXiv:2203.02167*, 2022.

[46] Chen Chen, Yufei Wang, Bing Li, and Kwok-Yan Lam. Knowledge is flat: A seq2seq generative framework for various knowledge graph completion. *arXiv preprint arXiv:2209.07299*, 2022.

[47] Chen Chen et al. Bridging structure and text for effective knowledge graph completion via conditional soft prompting. *arXiv preprint arXiv:2307.01709*, 2023.

[48] Yanbin Wei et al. Kicgpt: Large language model with knowledge in context for kg completion. *arXiv preprint arXiv:2402.02389*, 2024.

[49] Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. *arXiv preprint arXiv:1902.10197*, 2019.

[50] Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. Quaternion knowledge graph embeddings. *Advances in neural information processing systems*, 32, 2019.

[51] Donghan Yu and Yiming Yang. Retrieval-enhanced generative model for large-scale knowledge graph completion. In *Proceedings of SIGIR*, pages 2334–2338, 2023.

[52] Ruobing Xie, Zhiyuan Liu, Jia Jia, Huanbo Luan, and Maosong Sun. Representation learning of knowledge graphs with entity descriptions. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.

[53] Marko Pavlović, Yuta Wang, and Xia Li. Expressive: A spatio-functional embedding for knowledge graph completion. In *International Conference on Learning Representations (ICLR)*, 2023.

[54] Hongbo Liu, Yifan Zhang, and Lixin Zhou. An inductive knowledge graph embedding via combination of subgraph and type information. *Scientific Reports*, 13(1):12345, 2023.

[55] Peng He, Jie Sun, and Rui Chen. A type-augmented knowledge graph embedding framework for knowledge graph completion. *Scientific Reports*, 13(2):23456, 2023.

[56] Bosung Kim, Taesuk Hong, Youngjoong Ko, and Jungyun Seo. Multi-task learning for knowledge graph completion with pre-trained language models. In *Proceedings of COLING*, pages 1737–1743, 2020.

[57] Weijie Liu et al. K-bert: Enabling language representation with knowledge graph. In *Proceedings of AAAI*, pages 2901–2908, 2020.

[58] Louis Clouatre, Philippe Trempe, Amal Zouaq, and Sarath Chandar. Mlmlm: Link prediction with mean likelihood masked language model. *arXiv preprint arXiv:2009.07058*, 2020.

[59] Xin Lv et al. Do pre-trained models benefit knowledge graph completion? a reliable evaluation and a reasonable approach. In *Proceedings of ACL*, 2022.

[60] Bo Wang et al. Structure-augmented text representation learning for efficient knowledge graph completion. In *Proceedings of WWW*, pages 1737–1748, 2021.

[61] Da Li et al. Multi-task pre-training language model for semantic network completion. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 22(11), 2023.

[62] Justin Lovelace et al. Robust knowledge graph completion with stacked convolutions and a student re-ranking network. In *Proceedings of ACL*, pages 1016–1029, 2021.

[63] Daniel Daza, Michael Cochez, and Paul Groth. Inductive entity representations from text via link prediction. In *Proceedings of WWW*, pages 798–808, 2021.

[64] Yichi Zhang, Zhuo Chen, Wen Zhang, and Huajun Chen. Making large language models perform better in knowledge graph completion. *arXiv preprint arXiv:2310.06671*, 2023.

[65] Russa Biswas, Radina Sofronova, Mehwish Alam, and Harald Sack. Contextual language models for knowledge graph completion. In *Proceedings of MLSMKG@ECML-PKDD*, pages 13–24, 2021.

[66] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the AAAI conference on artificial intelligence*, volume 28, 2014.

[67] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.

[68] Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. Knowledge graph embedding via dynamic mapping matrix. In *Proceedings of the 53rd annual meeting of the association for computational linguistics and the 7th international joint conference on natural language processing (volume 1: Long papers)*, pages 687–696, 2015.

[69] Dai Quoc Nguyen, Dat Quoc Nguyen, Tu Dinh Nguyen, and Dinh Phung. A convolutional neural network-based model for knowledge base completion and its application to search personalization. *Semantic Web*, 10(5):947–960, 2019.

[70] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *The semantic web: 15th international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings 15*, pages 593–607. Springer, 2018.

[71] Liwei Cai and William Yang Wang. Kbgan: Adversarial learning for knowledge graph embeddings. *arXiv preprint arXiv:1711.04071*, 2017.

[72] Haonan Lu and Hailin Hu. Dense: An enhanced non-abelian group representation for knowledge graph embedding. *arXiv preprint arXiv:2008.04548*, 2020.

[73] Yanhui Peng and Jing Zhang. Lineare: Simple but powerful knowledge graph embedding for link prediction. In *2020 IEEE international conference on data mining (ICDM)*, pages 422–431. IEEE, 2020.

[74] Zhanqiu Zhang, Jianyu Cai, and Jie Wang. Duality-induced regularizer for tensor factorization based knowledge graph completion. *Advances in Neural Information Processing Systems*, 33: 21604–21615, 2020.

[75] Yushi Bai, Zhitao Ying, Hongyu Ren, and Jure Leskovec. Modeling heterogeneous hierarchies with relation-specific hyperbolic cones. *Advances in Neural Information Processing Systems*, 34:12316–12327, 2021.

[76] Tengwei Song, Jie Luo, and Lei Huang. Rot-pro: Modeling transitivity by projection in knowledge graph embedding. *Advances in Neural Information Processing Systems*, 34:24695–24706, 2021.

[77] Bonggeun Choi, Daesik Jang, and Youngjoong Ko. Mem-kgc: Masked entity model for knowledge graph completion with pre-trained language model. *IEEE Access*, 9:132025–132032, 2021.

[78] Apoorv Saxena, Adrian Kochsiek, and Rainer Gemulla. Sequence-to-sequence knowledge graph completion and question answering. *arXiv preprint arXiv:2203.10321*, 2022.

[79] Bonggeun Choi and Youngjoong Ko. Knowledge graph extension with a pre-trained language model via unified learning method. *Knowledge-Based Systems*, 262:110245, 2023.

[80] Xintao Wang, Qianyu He, Jiaqing Liang, and Yanghua Xiao. Language models as knowledge embeddings. *arXiv preprint arXiv:2206.12617*, 2022.

[81] Xin Xie et al. From discrimination to generation: Knowledge graph completion with generative transformer. In *Proceedings of WWW Companion*, pages 162–165, 2022.

# NeurIPS Paper Checklist

1. **Claims**

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

   Answer: [Yes]

   Justification: Yes, the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope.

   Guidelines:

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

2. **Limitations**

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

   Answer: [Yes]

   Justification: Yes, we discuss the limitations in Appendix F.

   Guidelines:

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

3. **Theory assumptions and proofs**

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

   Answer: [NA]

Justification: This paper does not include theoretical results.

Guidelines:

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

4. **Experimental result reproducibility**

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

Answer: [Yes]

Justification: We propose a novel KRL method in Section 3. We provide more details of the module implementation in Appendix C. Our source code is accessible in Abstract.

Guidelines:

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

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

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Yes, our source code is accessible in Abstract. All datasets used for evaluation are open-sourced.

Guidelines:

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

6. **Experimental setting/details**

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

Answer: [Yes]

Justification: Yes, we provide detailed descriptions of the experimental and evaluation settings in Section 4 and Appendix E.

Guidelines:

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

7. **Experiment statistical significance**

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

Answer: [No]

Justification: Although we do not report error bars, we provide source code in the Abstract.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

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

   Answer: [Yes]

   Justification: Yes, we report the details of experiments compute resources in Appendix E.

   Guidelines:
   - The answer NA means that the paper does not include experiments.
   - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
   - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
   - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

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

   Answer: [Yes]

   Justification: Yes, we have reviewed the NeurIPS Code of Ethics and ensured full compliance throughout our research process.

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

10. **Broader impacts**

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

    Answer: [Yes]

    Justification: Yes, we discuss the potential positive societal impacts and negative societal impacts in Appendix G.

    Guidelines:
    - The answer NA means that there is no societal impact of the work performed.
    - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
    - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

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

Answer: [NA]

Justification: This paper poses no such risks.

Guidelines:

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

12. **Licenses for existing assets**

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

Answer: [Yes]

Justification: Yes, we cite the original papers that produced the code package or dataset in Section 4.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

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

    Answer: [Yes]

    Justification: Yes, we communicate the details of the code as part of our submission. Our source code is anonymous.

    Guidelines:

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

14. **Crowdsourcing and research with human subjects**

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

    Answer: [Yes]

    Justification: This paper does not involve crowdsourcing nor research with human subjects.

    Guidelines:

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

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

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

    Answer: [NA]

    Justification: This paper does not involve crowdsourcing nor research with human subjects.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
    - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
    - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

    Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

    Answer: [Yes]

    Justification: LLMs were employed to support the understanding and interpretation of complex technical concepts.

    Guidelines:

    - The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
    - Please refer to our LLM policy (`https://neurips.cc/Conferences/2025/LLM`) for what should or should not be described.

# A Notations

| Notation | Explanation |
| --- | --- |
| $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T})$ | A knowledge graph consisting of entities $\mathcal{E}$, relations $\mathcal{R}$, and triples $\mathcal{T}$. |
| $\mathcal{E}$ | Set of entities in the knowledge graph. |
| $\mathcal{R}$ | Set of relations in the knowledge graph. |
| $\mathcal{T} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ | Set of relational triples $(e_i, r, e_j)$. |
| $(e_i, r, e_j)$ | A triple linking head entity $e_i$ to tail entity $e_j$ by relation $r$. |
| $\mathcal{V}$ | Textual vocabulary extracted from large-scale corpora. |
| $d$ | Shared embedding dimension for all entities and tokens. |
| $e \in \mathcal{E}$ | An individual entity. |
| $v \in \mathcal{V}$ | An individual vocabulary token. |
| $\mathbf{x} = (x_1, \ldots, x_N)$ | Tokenised input sequence of length $N$. |
| $N$ | Sequence length (number of tokens). |
| $\Theta$ | All trainable model parameters. |
| $\mathbf{H}^{(0)}$ | Initial embedding matrix for the input sequence. |
| $L$ | Number of stacked transformer layers. |
| $h$ | Number of attention heads in multi-head attention. |
| $d_h$ | Dimensionality of each attention head ($hd_h = d$). |
| $\mathbf{H}^{(\ell)} \in \mathbb{R}^{N \times d}$ | Hidden states at transformer layer $\ell$. |
| $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ | Query, key and value matrices in self-attention. |
| $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ | Projection parameters mapping hidden states to $\mathbf{Q}, \mathbf{K}, \mathbf{V}$. |
| $\mathbf{M}$ | Standard attention mask (causal or fully bidirectional). |
| $\mathbf{M}^{\text{BKA}}$ | Graph-aware bidirectional knowledge-attention mask (Eq. 2). |
| $\text{Attention}(\cdot)$ | Scaled-dot-product attention with mask $\mathbf{M}^{\text{BKA}}$. |
| $\tilde{\mathbf{H}}^{(\ell)}$ | Concatenated multi-head attention output at layer $\ell$. |
| $\mathbf{W}_O$ | Output projection matrix after concatenating heads. |
| $\bar{\mathbf{H}}^{(\ell)}$ | Hidden states after residual connection and layer normalisation. |
| $\text{FFN}(\cdot)$ | Position-wise feed-forward network. |
| $\mathbf{H}^{(\ell+1)}$ | Output hidden states of layer $\ell$ after FFN. |
| $\mathcal{M}$ | Index set of masked positions in knowledge-masked prediction. |
| $\mathbf{x}^m$ | Input sequence with tokens at $\mathcal{M}$ replaced by $\langle \text{mask} \rangle$. |
| $\gamma = \frac{|\mathcal{M}|}{N}$ | Masking ratio. |
| $\mathbf{h}_i$ | Final hidden representation at position $i$. |
| $p_\Theta(x_i \mid \mathbf{h}_i)$ | Predicted distribution over tokens/entities for position $i$. |
| $\mathbf{W}_P, \mathbf{b}_P$ | Output projection parameters for masked prediction. |
| $\mathcal{L}_{\text{KMP}}$ | Knowledge-masked prediction loss (Eq. 10). |
| $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}$ | Two independently corrupted views of the same sample for contrastive learning. |
| $\mathbf{H}^{(L)}(\mathbf{x}^{(k)})$ | Final-layer hidden states for view $k \in \{1, 2\}$. |
| $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}$ | Pooled sequence-level embeddings of the two views. |
| $\text{Pool}(\cdot)$ | Aggregation function (e.g. mean-pooling or [CLS] token). |
| $B$ | Minibatch size for contrastive learning. |
| $\text{sim}(\mathbf{u}, \mathbf{v})$ | Cosine similarity between vectors $\mathbf{u}$ and $\mathbf{v}$. |
| $\mathcal{L}_{\text{CGSA}}$ | Contrastive graph semantic aggregation loss (Eq. 12). |
| $\tau$ | Temperature hyperparameter in the contrastive loss. |

# B Related Work

## B.1 Knowledge Graph Embeddings

Recent KGE research has moved beyond early translational families to more expressive geometries and inductive designs that strive for stronger reasoning and generalisation. Nonetheless, purely structural models still struggle to capture contextual semantics and to operate in open–world settings. ExpressivE [53] embeds entities as points and relations as parallelograms in a learnable Euclidean sub-space, enabling the model to reproduce relation patterns such as hierarchy or composition with elegant geometric operations. TGraiL [54] combines sub-graph structure and type information to

support inductive link prediction where novel entities appear only at test time, while TaKE [55] injects latent type vectors into any backbone KGE model and devises a type-aware negative-sampling strategy to stabilise training. Although the above methods markedly improve structural fidelity, they inherit two long-standing weaknesses: (i) scarce usage of contextual semantics (e. g., textual descriptions) means that zero-shot prediction remains brittle; (ii) most models must be retrained whenever a new graph is introduced, hampering scalability. KG-BiLM retains the high information density of geometric encoders yet enriches the embeddings with bidirectional linguistic cues supplied by a large language model (LLM). In doing so, the framework alleviates the semantic sparsity and domain-adaptation barriers observed above.

These studies highlight the need for more flexible architectures that can effectively handle heterogeneous downstream scenarios that increasingly extends beyond entity-level embedding precision to more generalizable and context-sensitive representations.

### B.2 Encoder-based KRL

Encoder-based KRL methods treat a knowledge-graph (KG) triple as natural language and rely on pre-trained language models (PLMs) to inject rich semantic cues that traditional translational or tensor-factorisation approaches cannot capture. Broadly, existing work can be grouped into (i) triple-based representation, (ii) translation-based representation, and (iii) independent representation, each optimising the interplay between structural and textual signals in a different way

**Triple-Based Representation**. Triple-based models linearise an entire triple $(h, r, t)$ into a single sentence and fine-tune an encoder to judge its plausibility. KG-BERT [27] pioneers this paradigm: the description of the head entity (HE), relation (R) and tail entity (TE) is concatenated as [!CLS] HE [!SEP] R [!SEP] TE [!SEP] and fed to BERT; the [CLS] embedding is scored with a sigmoid-activated classifier. Subsequent work focuses on alleviating two key bottlenecks—insufficient relational modelling and lexical confusion between candidates. MTL-KGC [56] introduces multi-task objectives (relation prediction and relevance ranking) that expose KG-specific inductive bias during fine-tuning and markedly improve Hits@K. K-BERT [57] injects triples from an external KG into the input tree with soft position embeddings, effectively fusing factual memories into the encoder without retraining the PLM from scratch. MLMLM [58] reformulates link prediction as masked-language modelling (MLM): entities are generated token-by-token, enabling open-vocabulary completion and better interpretability. PKGC [59] and CSProm-KG [47] leverage prompt engineering; PKGC converts triples into cloze-style natural language while CSProm-KG learns conditional soft prompts that adapt to graph structure, reducing the parameter-update footprint and providing stronger zero-shot generalisation.

**Translation-Based Representation**. Translation-style encoders mimic the geometric intuition of TransE while retaining PLM expressiveness. A representative architecture is StAR [60], which encodes the $(h, r)$ pair and $t$ separately using a Siamese BERT and enforces that the pooled vectors satisfy $\|u - v\|_2$ margin constraints. The InfoNCE-like contrastive term encourages uniform spacing of negative samples in the joint space. Efficiency and negative-sampling quality are two recurring themes: SimKGC [45] introduces in-batch, pre-batch, and self negatives inside a bi-encoder, scaling contrastive learning to millions of triples without quadratic complexity. LP-BERT [61] generalises this idea with a masked entity–relation modelling task, facilitating inductive reasoning over previously unseen entities and relations. These works demonstrate that aligning head-relation and tail embeddings in a shared metric space remains competitive when enriched with contextual semantics.

**Independent Representation**. Independent or component-wise encoding disentangles the triple into three sequences, giving the model maximal flexibility to recombine learned embeddings. KEPLER [28] attaches a special token to each entity description, maps relations to trainable vectors, and jointly optimises a translational distance loss and Wikipedia-based MLM. This dual training scheme mitigates the frequency imbalance problem of long-tail entities. Two notable extensions illustrate the breadth of this design space: BERT-ResNet [62] stacks residual CNN blocks on top of contextualised embeddings to capture higher-order neighbourhood signals, boosting robustness on sparsely connected biomedical KGs. BLP [63] targets inductive link prediction by dynamically constructing mini-graphs around unseen entities and encoding them through adaptive graph convolutions plus PLM features.

However, such approaches remain narrow in focus, relying on local context windows that limit coverage of broader semantic relationships across the KG, ultimately constraining their ability to fully integrate structure and textual cues.

## B.3  Decoder-based KRL

Decoder-based KRL departs from encoder-only paradigms by exploiting the autoregressive decoders of large language models (LLMs) such as LLaMA-2 [12] and GPT-4 [30]. Because the decoder has been pre-exposed to vast world knowledge, these methods can generate or evaluate textual surrogates of triples with little or no parameter tuning. We group the literature into (i) description-generation, (ii) prompt-engineering, and (iii) structural fine-tuning.

**Description-Generation**. Low-resource entities suffer from sparse textual context, which traditional encoders rely on for semantics; decoder models therefore hallucinate rich descriptions that downstream KGC modules can ingest. Contextualisation Distillation (CD) [34] converts a raw triple $(h, r, t)$ into a natural-language prompt (e.g., "Describe the relationship between $h$ and $t$ given that $r$ holds."). An LLM produces a paragraph-length context $c$. Two auxiliary tasks—masked reconstruction of $c$ and generation of $c$ from the original triple—train a compact student model; weighting terms $\alpha$ and $\beta$ in the final loss encourage balanced learning of structure and language. CP-KGC [35] follows a similar philosophy but emphasises constrained prompts that prevent factual drift. The authors design slot-filled templates whose lexical space is restricted by ontology types, dramatically reducing hallucinated entities. By regenerating or extending existing descriptions, CP-KGC cuts long-tail error rates by 20 percents relative to free-form generation. These results highlight that what to ask the decoder is as important as how powerful the decoder is.

**Prompt-Engineering**. Prompt-engineering reframes link prediction or triple classification as a zero- or few-shot QA problem. KG-LLM [33] converts a triple into a yes/no question ("Is it true that Paris is-the-capital-of France?") and lets an instruction-tuned LLaMA decide plausibility. The scalar probability returned by the decoder acts as the score $s(h, r, t)$. With a modest 600 example instruction-tuning set, KG-LLM reaches TransE-era performance while requiring no negative sampling and no task-specific classifier. Subsequent work demonstrates that the same prompt can be morphed into an open-ended question to produce tail candidates directly, turning the decoder into a generative search engine. KICGPT [48] pushes efficiency further via in-context learning (ICL). At inference time, a handful of retrieved neighbour triples are embedded in the prompt as knowledge shots; no gradient updates are made. A structure-aware retriever prioritises paths that terminate with long-tail entities, mitigating the bias toward head-entities stored in the LLM's pre-training corpus. Compared with conventional fine-tuning, KICGPT reduces GPU hours by two orders of magnitude and still improves MRR on NELL-One by 6 percents. The same idea scales to multilingual scenarios—KG-LLaMA and KG-ChatGLM augment prompts with language tags and achieve reliable cross-lingual transfer without any machine-translation pipeline. These results confirm that the decoder's built-in language universality can be channelled toward structured reasoning with minimal engineering.

**Structural Fine-Tuning**. While pure prompting leverages implicit knowledge, structural fine-tuning injects explicit KG embeddings into the decoder. KoPA [64] first pre-trains TransE-style vectors $h, r, t$ and then transforms them into virtual prefix tokens via a lightweight adapter $P(\dot)$. The resulting sequence is fed to the LLM, which learns to condition its generation on both text and structure. Crucially, only the adapter parameters are updated; the frozen LLaMA backbone retains its linguistic fluency. On FB15k-237, KoPA boosts Hits@1 by 3 percents over adapter-free fine-tuning and narrows the gap between parameter-efficient tuning and full-model updates. KG-GPT2 [65] offers an earlier but influential proof-of-concept: each triple is linearised as a sentence ("head [SEP] relation [SEP] tail") and GPT-2 computes its perplexity as a plausibility score. Subsequent variants append neighbour triples as additional context tokens, enabling explainable chain-of-thought style completion. Finally, retrieval-augmented generators such as RESKGC [51] hybridise all three themes. A BM25 retriever first pulls textual evidence; GPT-3.5 then generates candidate triples, which are re-ranked by structural filters. Although RESKGC was developed for web-scale graphs, its modular pipeline hints at future decoder systems that couple fast symbolic indices with slow but powerful reasoners.

These studies provide important insights into the synergy of generative modeling and structured data. Nevertheless, the methods mostly output textual responses without a dedicated mechanism to retain or highlight the structural integrity of knowledge graphs, thus limiting their capacity for robust graph-based inference.

---

**Algorithm 1:** Bidirectional Knowledge Attention (BKA)

---

**Input:** Token sequence $\mathbf{x} = (x_1, \ldots, x_N)$; knowledge graph $\mathcal{G} = (\mathcal{E}, \mathcal{R}, \mathcal{T})$; hop threshold $h$;
  model parameters $\{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O\}$

**Output:** Contextually and structurally enriched representations $\mathbf{H}^{(L)}$

---

Initialize hidden states $\mathbf{H}^{(0)} \leftarrow \text{Embed}(\mathbf{x})$;
**for** $\ell \leftarrow 0$ **to** $L - 1$ **do**
 **for** $i \leftarrow 1$ **to** $N$ **do**
  $\mathbf{q}_i \leftarrow \mathbf{H}_i^{(\ell)} \mathbf{W}_Q$;
  **for** $j \leftarrow 1$ **to** $N$ **do**
   $\mathbf{k}_j \leftarrow \mathbf{H}_j^{(\ell)} \mathbf{W}_K$;
   **if** $\text{CanInteract}(x_i, x_j, \mathcal{G}, h)$ **then**
    $\mathbf{M}_{i,j}^{\text{BKA}} \leftarrow 0$;
   **else**
    $\mathbf{M}_{i,j}^{\text{BKA}} \leftarrow -\infty$;

  $\boldsymbol{\alpha}_i \leftarrow \text{softmax}\left(\frac{\mathbf{q}_i \mathbf{K}^{\top}}{\sqrt{d_h}} + \mathbf{M}_{i,:}^{\text{BKA}}\right)$;
  $\mathbf{V} \leftarrow \mathbf{H}^{(\ell)} \mathbf{W}_V$;
  $\tilde{\mathbf{h}}_i \leftarrow \boldsymbol{\alpha}_i \mathbf{V}$;
 $\tilde{\mathbf{H}}^{(\ell)} \leftarrow \text{Concat}(\tilde{\mathbf{h}}_1, \ldots, \tilde{\mathbf{h}}_N) \mathbf{W}_O$;
 $\bar{\mathbf{H}}^{(\ell)} \leftarrow \text{LayerNorm}(\mathbf{H}^{(\ell)} + \tilde{\mathbf{H}}^{(\ell)})$;
 $\mathbf{H}^{(\ell+1)} \leftarrow \text{LayerNorm}(\bar{\mathbf{H}}^{(\ell)} + \text{FFN}(\bar{\mathbf{H}}^{(\ell)}))$;
**return** $\mathbf{H}^{(L)}$;

---

## C  Details of KG-BiLM Modules

### C.1  Bidirectional Knowledge Attention

The BKA algorithm constitutes the connective tissue that fuses a Transformer's bidirectional reasoning capacity with the multi-hop relational structure inherent in a KG. At its core, BKA generalises the classical attention mechanism by augmenting the binary attention mask with graph-aware semantics, thereby allowing each token to attend not only to every other token in its textual context but also to any entity token that is reachable by a bounded number of relational hops. The resulting attention topology is simultaneously sequence-complete and graph-selective, which contrasts starkly with the hard uni-directionality of decoder-only language models and the structure-agnostic bidirectionality of pure encoders.

#### C.1.1  Design Rationale

Traditional KGE pipelines frequently rely on translation-based objectives that view relations as static vectors in $\mathbb{R}^d$ and treat textual aliases as exogenous after-thoughts. Conversely, text-centric LLMs are exceptional at capturing lexical regularities but lack explicit exposure to the symbolic constraints of a KG—symmetry, transitivity, inverse relations, and so forth. BKA is conceived as a minimal yet expressive modification that marries these two worlds. By introducing the function CanInteract, the algorithm decides, for every ordered pair of positions $(i, j)$, whether their embeddings should participate in the same attention subspace. The decision is contingent on two orthogonal axes: temporal orientation (past, present, future tokens) and structural connectivity (graph adjacency up to $h$ hops).

In effect, BKA replaces the Transformer's conventional triangular (causal) mask or fully populated (encoder) mask with a knowledge-adaptive mask. When $h = 0$ and CanInteract tests only the inequality $j \leq i$, BKA reduces to a standard causal decoder. When $h = \infty$ and structural checks are disabled, BKA degenerates to a vanilla bidirectional BERT-style encoder. This flexible continuum allows practitioners to choose a sweet-spot where syntactic coherence and graph coherence are mutually reinforced.

### C.1.2 Algorithmic Steps

**Initial Embedding**. The first line in the pseudocode instantiates $\mathbf{H}^{(0)}$ by concatenating standard token embeddings with position encodings and, crucially, KG-entity type embeddings. Unlike typical LM practice, entity IDs are not merely treated as sub-word strings; rather, they possess dedicated vectors that are updated jointly with lexical embeddings.

**Graph-Aware Mask Construction**. Within each layer $\ell$, nested loops traverse token positions. For every pair $(i, j)$ the routine computes $\mathbf{M}_{i,j}^{\mathrm{BKA}}$. The helper CanInteract returns true if at least one of three conditions is satisfied: (1) $|i - j| \leq \delta$, where $\delta$ is a small local window to preserve micro-syntax; (2) $x_i$ and $x_j$ are verbal tokens and the model operates in a language-model mode; (3) their parent entities in the KG are connected by a path not exceeding $h$ hops. The final mask thus embodies both text contiguity and relational contiguity.

**Attention Weight Computation**. Once the mask matrix has been filled, the algorithm computes the scaled dot-product attention for each query. Because the mask contains $-\infty$ for disallowed pairs, the corresponding softmax weights vanish, enforcing hard structural sparsity. While the line-wise loops are explicit for clarity, real implementations vectorise all positions and heads, thereby preserving the $O(N^2)$ theoretical complexity of standard attention.

**Residual Path and Feed-Forward Layers**. After multi-head projection through $\mathbf{W}_O$, a residual connection and layer normalisation are applied. The FFN applies two linear maps with a gating non-linearity, typically GELU, and a residual link—a design that is empirically crucial for stabilising deep stacks.

### C.1.3 Theoretical Properties

**Symmetry Preservation**. Because the mask allows $i$ to see $j$ and vice versa whenever CanInteract is true, relational symmetry (e.g., "married to") can be internalised in a single pass rather than via post-hoc rule injection.

**Multi-Hop Relational Reasoning**. If the hop threshold $h > 1$, BKA implicitly permits the propagation of information along length-$h$ paths. The self-attention coefficients can be shown to approximate a truncated power series of the KG adjacency matrix, thereby endowing the representations with spectral properties reminiscent of Graph Neural Networks.

**Expressivity under Causal Decoding**. Despite affording bidirectional flows, BKA does not sabotage the autoregressive nature required for text generation. During inference, tokens at position $i$ only depend on already generated tokens plus KG neighbours, both of which are known at decoding time. Thus, BKA is compatible with beam search and other sampling schemes.

## C.2 Knowledge-Masked Prediction

KMP extends the masked-language modelling paradigm by integrating explicit structural priors from the knowledge graph and, crucially, by shifting the prediction locus from the masked position to its immediate predecessor. This seemingly minor alteration endows KG-BiLM with a potent capacity for causal reasoning under partial observability, thereby promoting robust zero-shot transfer.

### C.2.1 Conceptual Foundations

Standard MLM tasks such as those employed in BERT require the model to reconstruct randomly hidden tokens given an unconstrained bidirectional view. While effective for textual semantics, that setup is less compatible with autoregressive decoding and overlooks the rich relational patterns that connect entity mentions across a document. KMP addresses both issues simultaneously:

**Shifted Prediction**. By consuming $\mathbf{h}i - 1$ instead of $\mathbf{h}i$, the model must rely on context that excludes the target token, thereby simulating conditions at generation time. This choice also compels the network to exploit forward context captured via BKA—without it, $\mathbf{h}_{i-1}$ would lack information about tokens $j > i - 1$.

**Structure-Aware Masking**. Mask selection is not purely random. Entity positions that share high centrality or critical relation types (e.g., is-a, part-of) are preferentially masked with higher probability,

---

**Algorithm 2:** Knowledge-Masked Prediction (KMP)

---

**Input:** Token sequence $\mathbf{x} = (x_1, \ldots, x_N)$; mask ratio $\gamma$; projection parameters $\{\mathbf{W}_P, \mathbf{b}_P\}$;
   BKA-augmented Transformer
**Output:** Masked-prediction loss $\mathcal{L}_{\mathrm{KMP}}$

---

Sample mask set $\mathcal{M} \subseteq \{1, \ldots, N\}$ with $|\mathcal{M}| = \gamma N$;
**foreach** $i \in \mathcal{M}$ **do**
  Replace $x_i$ with $\langle \mathrm{mask} \rangle$ to obtain $\mathbf{x}^m$
$\mathbf{H}^{(L)} \leftarrow \mathrm{BKA\_Forward}(\mathbf{x}^m)$;
$\mathcal{L}_{\mathrm{KMP}} \leftarrow 0$;
**foreach** $i \in \mathcal{M}$ **do**
  $\mathbf{h}_{i-1} \leftarrow \mathbf{H}^{(L)}_{i-1}$;
  $\boldsymbol{\ell}_i \leftarrow \mathbf{W}_P \mathbf{h}_{i-1} + \mathbf{b}_P$;
  $p_\Theta(x_i | \mathbf{h}_{i-1}) \leftarrow \mathrm{softmax}(\boldsymbol{\ell}_i)$;
  $\mathcal{L}_{\mathrm{KMP}} += -\log p_\Theta(x_i | \mathbf{h}_{i-1})$;
$\mathcal{L}_{\mathrm{KMP}} \leftarrow \frac{1}{|\mathcal{M}|} \mathcal{L}_{\mathrm{KMP}}$;
**return** $\mathcal{L}_{\mathrm{KMP}}$;

---

calibrated by an importance score derived from PageRank and relation entropy. This targeted masking encourages the model to learn structural dependencies that matter for downstream link prediction.

### C.2.2 Algorithmic Steps

**Mask Construction**. A binary mask vector is sampled such that exactly $\gamma N$ positions are designated for replacement. For non-entity tokens, a uniform distribution suffices; for entity tokens, a weighted sampling scheme ensures that high-degree nodes are systematically challenged.

**Encoding via BKA**. The masked input $\mathbf{x}^m$ is fed through the BKA stack, which already embodies graph-informed attention patterns. Therefore, even though certain entity IDs have been replaced by an opaque mask token, their neighbours can still convey hints about their identity through multi-hop edges.

**Cross-Entropy Accumulation**. For each masked index $i$, logits are produced from $\mathbf{h}i-1$ and normalised via softmax over the joint vocabulary $\mathcal{V} \cup \mathcal{E}$. Crucially, the inclusion of entity symbols means the classifier is solving a mixed-type prediction problem. The average negative log-likelihood across all masked positions yields $\mathcal{L}$KMP.

**Backpropagation and Parameter Update**. Because the loss is defined at the level of the final projection, gradients propagate through both projection parameters and the entire BKA stack, refining lexical embeddings, entity embeddings, and relation embeddings in concert.

### C.3 Contrastive Graph Semantic Aggregation

Contrastive learning has emerged as a powerful unsupervised paradigm for representation learning, particularly in vision and speech domains. Its core premise is to pull together embeddings that originate from different views of the same underlying sample while pushing apart embeddings from distinct samples. In the context of knowledge-graph language models, however, naive contrastive approaches face two entrenched challenges. First, entities and relations seldom appear in isolation; their semantics are entangled with multi-hop relational paths whose topological signatures must be respected. Second, textual paraphrases and graph augmentations induce heterogeneous perturbations, making it non-trivial to judge when two views are "positive."

The proposed CGSA addresses these issues by marrying graph topology–aware augmentation with a Transformer encoder that is already enriched by BKA and trained under KMP. CGSA sits on top of this encoder and imposes a clustering prior on entire sub-graph representations, ensuring that embeddings remain discriminative despite repeated dropout or textual paraphrases. Conceptually, CGSA plays the role of a semantic pressure valve: it prevents the high-capacity encoder from collapsing entity neighbourhoods into an overly tangled representation space.

**Algorithm 3:** Contrastive Graph Semantic Aggregation (CGSA)

---

**Input:** Minibatch of subgraphs or text snippets $\mathcal{S} = \{s_1, \ldots, s_B\}$; temperature $\tau$; hop threshold $h$; BKA + KMP encoder $f_\Theta$; pooling operator $\text{Pool}(\cdot)$
**Output:** Contrastive loss $\mathcal{L}_{\text{CGSA}}$

---

$\mathcal{L}_{\text{CGSA}} \leftarrow 0$;
**foreach** $s_k \in \mathcal{S}$ **do**
    Generate two stochastic views by augmentation;
    $\mathbf{x}_k^{(1)} \leftarrow \text{Augment}(s_k, h)$;
    $\mathbf{x}_k^{(2)} \leftarrow \text{Augment}(s_k, h)$;
    $\mathbf{H}_k^{(L,1)} \leftarrow f_\Theta(\mathbf{x}_k^{(1)})$;
    $\mathbf{H}_k^{(L,2)} \leftarrow f_\Theta(\mathbf{x}_k^{(2)})$;
    $\mathbf{z}_k^{(1)} \leftarrow \text{Pool}(\mathbf{H}_k^{(L,1)})$;
    $\mathbf{z}_k^{(2)} \leftarrow \text{Pool}(\mathbf{H}_k^{(L,2)})$;
    $\mathbf{z}_k^{(1)} \leftarrow \mathbf{z}_k^{(1)} / \|\mathbf{z}_k^{(1)}\|$;
    $\mathbf{z}_k^{(2)} \leftarrow \mathbf{z}_k^{(2)} / \|\mathbf{z}_k^{(2)}\|$;
**for** $k \leftarrow 1$ **to** $B$ **do**
    $\text{pos} \leftarrow \exp\big(\text{sim}(\mathbf{z}_k^{(1)}, \mathbf{z}_k^{(2)})/\tau\big)$;
    $\text{den} \leftarrow 0$;
    **for** $\ell \leftarrow 1$ **to** $B$ **do**
        $\text{den} \mathrel{+}= \exp\big(\text{sim}(\mathbf{z}_k^{(1)}, \mathbf{z}_\ell^{(2)})/\tau\big)$;
    $\mathcal{L}_{\text{CGSA}} \mathrel{-}= \log\big(\text{pos}/\text{den}\big)$;
$\mathcal{L}_{\text{CGSA}} \leftarrow \mathcal{L}_{\text{CGSA}}/B$;
**return** $\mathcal{L}_{\text{CGSA}}$;

---

### C.3.1 Algorithmic Steps

Input to CGSA is a mini-batch of $B$ samples, each being either a textual snippet or a KG sub-graph. For clarity, we call each sample $s_k$ a semantic unit. The algorithm proceeds in two macro phases: (i) stochastic view generation and (ii) InfoNCE loss computation.

**Stochastic View Generation**. For every semantic unit $s_k$, CGSA calls an augmentation routine twice, producing $\mathbf{x}_k^{(1)}$ and $\mathbf{x}_k^{(2)}$. The augmentation operator obeys two design principles:

Graph Reachability: A hop threshold $h$ limits the radius within which entities may be dropped or substituted. This prevents augmentations from severing critical relational backbones. Textual Paraphrasing: For purely textual spans, augmentation may include synonym replacement, span deletion, or entity order shuffling, with the constraint that entity co-occurrence statistics are preserved up to the second moment.

Each augmented view is then passed through the encoder $f_\Theta$. Because $f_\Theta$ already employs BKA masks, it seamlessly integrates both local lexical features and multi-hop structural cues. The encoder returns token-level hidden states $\mathbf{H}_k^{(L,\cdot)}$, which are pooled into fixed-length vectors by a user-selected operator (e.g., mean pooling, max pooling, or [CLS] token extraction). Normalisation to unit length transforms the embeddings into points on the unit hypersphere—a prerequisite for cosine similarity to coincide with dot product.

**InfoNCE Loss Computation**. For each anchor index $k$, the algorithm forms a positive pair $(\mathbf{z}_k^{(1)}, \mathbf{z}_k^{(2)})$ and computes their scaled similarity, labelled pos. It then sums exponentiated similarities between the anchor $\mathbf{z}^{(1)}k$ and all second-view embeddings $\mathbf{z}^{(2)}\ell$ across the mini-batch, accumulating in den. The resulting expression implements the InfoNCE objective, which can be shown to maximise a lower bound on mutual information between the two views. Averaging over the batch yields the final loss $\mathcal{L}_{\text{CGSA}}$, later combined with BKA + KMP losses during multi-task training.

# D  Detailed Results

In this section, we provide the full experimental results for all baseline models on the evaluated datasets. Table 7 reports the link-prediction metrics on WN18RR and FB15k-237, while Table 8 presents the results on Wikidata5M and FB15k-237N.

Table 7: Link prediction metrics on WN18RR and FB15k-237 datasets

| Model | WN18RR | | | | | FB15k-237 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | MR | MRR | Hits@1 | Hits@3 | Hits@10 | MR | MRR | Hits@1 | Hits@3 | Hits@10 |
| TransE[14] | 2300 | 24.3 | 4.3 | 44.1 | 53.2 | 223 | 27.9 | 19.8 | 37.6 | 47.4 |
| TransH[66] | 2524 | – | – | – | 50.3 | 255 | – | – | – | 48.6 |
| DistMult[15] | 3704 | 44.4 | 41.2 | 47.0 | 50.4 | 411 | 28.1 | 19.9 | 30.1 | 44.6 |
| TransR[67] | 3166 | – | – | – | 50.7 | 237 | – | – | – | 51.1 |
| TransD[68] | 2768 | – | – | – | 50.7 | 246 | – | – | – | 48.4 |
| ComplEx[37] | 3921 | 44.9 | 40.9 | 46.9 | 53.0 | 508 | 27.8 | 19.4 | 29.7 | 45.0 |
| ConvE[38] | 4464 | 45.6 | 41.9 | 47.0 | 53.1 | 245 | 31.2 | 22.5 | 34.1 | 49.7 |
| ConvKB[69] | 2554 | 24.9 | – | – | 52.5 | 257 | 24.3 | – | – | 51.7 |
| R-GCN[70] | 6700 | 12.3 | 8.0 | 13.7 | 20.7 | 600 | 16.4 | 10.0 | 18.1 | 41.7 |
| KBGAN[71] | – | 21.5 | – | – | 48.1 | – | 27.7 | – | – | 45.8 |
| TuckER[39] | – | 47.0 | 44.3 | 48.2 | 52.6 | – | 35.8 | 26.6 | 39.4 | 54.4 |
| DensE[72] | 3052 | 49.1 | 44.3 | 50.8 | 57.9 | 169 | 34.9 | 25.6 | 38.4 | 53.5 |
| LineaRE[73] | 1644 | 49.5 | 45.3 | 50.9 | 57.8 | 155 | 35.7 | 26.4 | 39.1 | 54.5 |
| RESCAL-DURA[74] | – | 49.8 | 45.5 | – | 57.7 | – | 36.8 | 27.6 | – | 55.0 |
| CompGCN[40] | – | 47.9 | 44.3 | 49.4 | 54.6 | – | 35.5 | 26.4 | 39.0 | 53.5 |
| ConE[75] | – | 49.6 | 45.3 | 51.5 | 57.9 | – | 34.5 | 24.7 | 38.1 | 54.0 |
| Rot-Pro[76] | – | 45.7 | 39.7 | 48.2 | 57.7 | – | 34.4 | 24.6 | 38.3 | 54.0 |
| QuatDE[41] | 1977 | 48.9 | 43.8 | 50.9 | 58.6 | **90** | 36.5 | 26.8 | <u>40.0</u> | <u>56.3</u> |
| NBFNet[42] | – | 55.1 | 49.7 | – | 66.6 | – | **41.5** | 32.1 | – | **59.9** |
| KG-BERT[27] | 97 | 21.6 | 4.1 | 30.2 | 52.4 | 153 | 23.7 | 16.9 | 26.0 | 42.7 |
| MTL-KGC[56] | 89 | 33.1 | 20.3 | 38.3 | 59.7 | 132 | 26.7 | 17.2 | 29.8 | 45.8 |
| Pretrain-KGE[43] | – | 48.8 | 43.7 | 50.9 | 58.6 | – | 35.0 | 25.0 | 38.4 | 55.4 |
| StAR[60] | <u>51</u> | 40.1 | 24.3 | 49.1 | 70.9 | 1117 | 29.6 | 20.5 | 32.2 | 48.2 |
| MEM-KGC[77] | – | 57.2 | 48.9 | 62.0 | 72.3 | – | 34.9 | 26.0 | 38.2 | 52.4 |
| LaSS[44] | **35** | – | – | – | 78.6 | <u>108</u> | – | – | – | 53.3 |
| SimKGC[45] | – | 66.7 | 58.8 | <u>72.1</u> | **80.5** | – | 33.6 | 24.9 | 36.2 | 51.1 |
| LP-BERT[61] | 92 | 48.2 | 34.3 | 56.3 | 75.2 | 154 | 31.0 | 22.3 | 33.6 | 49.0 |
| KGT5[78] | – | 54.2 | 50.7 | – | 60.7 | – | 34.3 | 25.2 | – | 37.7 |
| OpenWorld KGC[79] | – | 55.7 | 47.5 | 60.4 | 70.4 | – | 34.6 | 25.3 | 38.1 | 53.1 |
| LMKE[80] | 79 | 61.9 | 52.3 | 67.1 | 78.9 | 141 | 30.6 | 21.8 | 33.1 | 48.4 |
| GenKGC[81] | – | – | 28.7 | 40.3 | 53.5 | – | – | 19.2 | 35.5 | 43.9 |
| KG-S2S[46] | – | 57.4 | 53.1 | 59.5 | 66.1 | – | 33.6 | 25.7 | 37.3 | 49.8 |
| kNN-KGE[29] | – | 57.9 | 52.5 | – | – | – | 28.0 | **37.3** | – | – |
| CSPromp-KG[47] | – | 57.5 | 52.2 | 59.6 | 67.8 | – | 35.8 | 26.9 | 39.3 | 53.8 |
| GPT-3.5 | – | – | 19.0 | – | – | – | – | 23.7 | – | – |
| CP-KGC[35] | – | <u>67.3</u> | <u>59.9</u> | <u>72.1</u> | <u>80.4</u> | – | 33.8 | 25.1 | 36.5 | 51.6 |
| KICGPT[48] | – | 56.4 | 47.8 | 61.2 | 67.7 | – | <u>41.2</u> | <u>32.7</u> | **44.8** | 55.4 |
| KG-BiLM(Ours) | 67 | **68.2** | **61.4** | **72.7** | **80.5** | 151 | 36.7 | 30.5 | 36.9 | 53.1 |

# E  Implementation Details

**Model Configuration** All experiments instantiate KG-BiLM with $L = 24$ transformer layers, model dimension $d = 1,024$, $h = 16$ attention heads (so $d_h = 64$ per head), and a feed-forward hidden size of 4,096. We set the maximum sequence length to 512 tokens, and represent both entity and vocabulary embeddings in the same $d$-dimensional space. The hop threshold for the Bidirectional Knowledge Attention mask is set to 2 (i.e. entities within two hops in $\mathcal{G}$ may attend to each other).

**Training Hyperparameters** We optimize with Adam ($\beta_1$=0.9, $\beta_2$=0.999, $\epsilon = 10^{-8}$) and a linear learning-rate schedule with 10k warm-up steps. The peak learning rate is $1 \times 10^{-4}$, decayed to zero over 200k total steps. We apply a weight decay of $1 \times 10^{-2}$ and gradient clipping at norm 1.0. Dropout of 0.1 is used in both attention and feed-forward sublayers. For KMP, the masking ratio $\gamma$ is 15

Table 8: Link prediction metrics on Wikidata5M and FB15k-237N datasets

| Model | Wikidata5M | | | | FB15k-237N | | | |
|---|---|---|---|---|---|---|---|---|
| | MRR | Hits@1 | Hits@3 | Hits@10 | MRR | Hits@1 | Hits@3 | Hits@10 |
| TransE[14] | 25.3 | 17.0 | 31.1 | 39.2 | 25.5 | 15.2 | 30.1 | 45.9 |
| DistMult[15] | 25.3 | 20.9 | 27.8 | 33.4 | 20.9 | 14.3 | 23.4 | 33.0 |
| ComplEx[37] | 30.8 | 25.5 | – | 39.8 | 24.9 | 18.0 | 27.6 | 38.0 |
| DKRL[52] | 16.0 | 12.0 | 18.1 | 22.9 | – | – | – | – |
| RoBERTa[24] | 0.1 | 0 | 0.1 | 0.3 | – | – | – | – |
| RotatE[49] | 29.0 | 23.4 | 32.2 | 39.0 | 27.9 | 17.7 | 32.0 | 48.1 |
| QuatE[50] | 27.6 | 22.7 | 30.1 | 35.9 | – | – | – | – |
| ConvE[38] | – | – | – | – | 27.3 | 19.2 | 30.5 | 42.9 |
| CompGCN[40] | – | – | – | – | 31.6 | 23.1 | 34.9 | 48.0 |
| KG-BERT[27] | – | – | – | – | 20.3 | 13.9 | 20.1 | 40.3 |
| MTL-KGC[56] | – | – | – | – | 24.1 | 16.0 | 28.4 | 43.0 |
| GenKGC[81] | – | – | – | – | – | 18.7 | 27.3 | 33.7 |
| KG-S2S[46] | – | – | – | – | 35.4 | 28.5 | 38.8 | 49.3 |
| KEPLER[28] | 21.0 | 17.3 | 22.4 | 27.7 | – | – | – | – |
| SimKGC[45] | 35.8 | 31.3 | 37.6 | 44.1 | – | – | – | – |
| KGT5[78] | 33.6 | 28.6 | 36.2 | 42.6 | – | – | – | – |
| CSPromp-KG[47] | 38.0 | 34.3 | 39.9 | <u>44.6</u> | 36.0 | 28.1 | 39.5 | 51.1 |
| ReSKGC[51] | <u>39.6</u> | <u>37.3</u> | <u>41.3</u> | 43.7 | – | – | – | – |
| CD[34] | – | – | – | – | <u>37.2</u> | <u>28.8</u> | <u>41.0</u> | <u>53.0</u> |
| **KG-BiLM(Ours)** | **40.3** | **39.7** | **43.0** | **45.2** | **37.8** | **29.3** | **42.1** | **54.6** |

**Contrastive Learning Settings** In CGSA, we sample mini-batches of $B = 256$ sub-graphs/snippets, generating two corrupted views each via independent dropout masks. We use mean-pooling over the final hidden states to obtain **z**-vectors, and set the temperature $\tau$ to 0.07.

**Hardware and Software Environment** All models are trained on a cluster of two NVIDIA H100 GPUs (80 GB HBM3 each) interconnected via NVLink. The software stack comprises PyTorch 2.0, CUDA 12.8, and NCCL 2.16 on Centos 9 Stream. The host features dual Intel Xeon Silver 4416 CPUs (80 cores total) and 512 GB RAM.

# F  Limitations

Although KG-BiLM delivers state-of-the-art results on three of the four benchmarks, its performance on the structure-only FB15k-237 dataset remains merely on par with recent transformer baselines and still trails the specialized path-based model NBFNet. The shortfall highlights two intrinsic limitations of our current design. (i) Relation-cardinality sensitivity. FB15k-237 contains more than twenty times as many distinct relations as WN18RR, yet offers no lexical clues. Under this setting, the Bidirectional Knowledge Attention module must rely exclusively on topological co-occurrence signals, and its ability to disambiguate semantically similar but label-distinct relations is attenuated. (ii) Semantic sparsity dependence. KG-BiLM's masking-recovery objective was tuned with the assumption that at least weak textual context would be available. When such cues are completely absent, the model still outperforms distance-based KGEs but cannot fully exploit its contrastive semantic aggregation, leading to diminished gains.

# G  Broader Impacts

The development of KG-BiLM—a hybrid architecture that unifies knowledge-graph structural representations with the semantic understanding capabilities of large language models—carries both promising benefits for society and important risks that must be carefully managed. In this section, we outline potential positive societal impacts alongside possible negative consequences, considering ethical, technical, and environmental dimensions.

**Enhanced Access to Domain Knowledge** By integrating rich graph structure with contextual semantics, KG-BiLM can serve as a powerful tool for making specialized information more accessible.

In domains such as healthcare, legal reasoning, or scientific research, practitioners often need to traverse complex networks of interrelated concepts. KG-BiLM's high-density embeddings and zero-shot capabilities can support more intuitive query interfaces, enabling non-experts to retrieve precise, contextually relevant information without extensive domain training.

**Improved Downstream Applications** The unified representations produced by KG-BiLM can boost the performance of a wide array of downstream tasks—from question answering and information extraction to recommendation systems and decision support. In education, for example, intelligent tutoring systems built on KG-BiLM could provide tailored explanations by jointly leveraging pedagogical ontologies and linguistic context. Likewise, in environmental monitoring, KG-BiLM could help integrate sensor data with domain taxonomies to surface early warnings of ecological disturbances.

**Facilitation of Interdisciplinary Research**. KG-BiLM's ability to align heterogeneous graph modalities with free-text semantics encourages cross-disciplinary collaboration. Researchers from different fields who maintain separate knowledge bases (e.g., biomedical ontologies and social science taxonomies) can benefit from a shared latent space that respects both structural constraints and nuanced textual insights. This could accelerate innovation at the intersections of AI, biology, and social policy.