

UniHR: Hierarchical Representation Learning for Unified Knowledge Graph Link Prediction

Anonymous ACL submission

Abstract

Beyond-triple fact representations including hyper-relational facts with auxiliary key-value pairs, temporal facts with additional timestamps, and nested facts implying relationships between facts, are gaining significant attention. However, existing link prediction models are usually designed for one specific type of facts, making it difficult to generalize to other types of facts. To overcome this limitation, we propose a **Unified Hierarchical Representation learning framework (UniHR)** for unified knowledge graph link prediction. It consists of a unified Hierarchical Data Representation (HiDR) module and a unified Hierarchical Structure Learning (HiSL) module as graph encoder. The HiDR module unifies hyper-relational KGs, temporal KGs, and nested factual KGs into triple-based representations. Then HiSL incorporates intra-fact and inter-fact message passing, focusing on enhancing the semantic information within individual facts and enriching the structural information between facts. Empirical results demonstrate the effectiveness of UniHR and highlight the strong potential of unified representations. Code and data are available at <https://anonymous.4open.science/r/UniHR-BDCB/>.

1 Introduction

Large-scale knowledge graphs (KGs) such as Freebase (Bollacker et al., 2008) and Wikidata (Vrandečić and Krötzsch, 2014) have been widely applied in many areas like question answering (Kaiser et al., 2021) and natural language processing (Anervaz et al., 2018). However, the presence of missing facts within these KGs inevitably limit their applications. Therefore, the link prediction task has been introduced to predict missing elements within factual data. Current link prediction methods mainly focus on the fact in the form of triple e.g., (*head entity*, *relation*, *tail entity*).

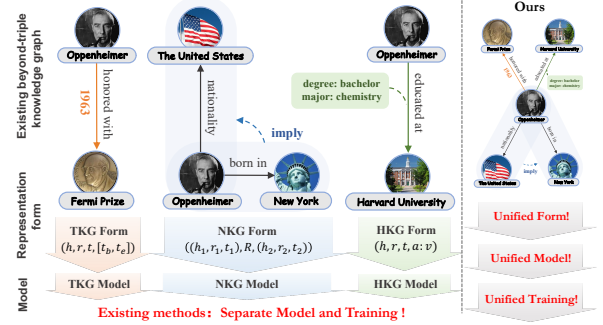


Figure 1: Comparison between the existing link prediction methods for beyond-triple KGs and our methods.

Despite the simplicity and unity of triple-based representation, it is difficult to adequately express complex facts, such as “*Oppenheimer is educated at Harvard University for a bachelor degree in chemistry*” shown in Figure 1. Therefore, existing researches (Wang et al., 2021; Xiong et al., 2024; Xu et al., 2019) contribute to focusing on semantically richer facts. Figure 1 illustrates three specific types of facts: hyper-relational fact ((*Oppenheimer*, *educated at*, *Harvard University*), *degree: bachelor*, *major: chemistry*), temporal fact (*Oppenheimer*, *honored with*, *Fermi Prize*, *1963*), nested fact ((*Oppenheimer*, *born in*, *New York*), *imply*, (*Oppenheimer*, *nationality*, *The United States*)). These forms of facts allow for expression of complex semantics and revelation of relationships between facts, extending beyond the triple-based representation. Thus in recent years, Hyper-relational KGs (HKG) (Chung et al., 2023), Temporal KGs (TKG) (Xu et al., 2023a), and Nested factual KGs (NKG) (Xiong et al., 2024) attract wide research interests.

Recent studies have demonstrated the effectiveness of various embedding strategies for these beyond-triple representations (Xiong et al., 2023). However, these methods are usually designed for specific representation forms, e.g., StarE (Galkin et al., 2020) customizes graph neural network to enhance structure information between hyper-

relational facts, For NKGs, BiVE (Chung and Whang, 2023) bridges semantics between two levels of facts through a simple MLP. Moreover, GeomE+ (Xu et al., 2023a) et al. TKG methods employ time-aware scoring functions to capture semantic information within individual temporal fact with timestamp. Although these methods perform well on specific types of facts, it is evident that such customized methods are difficult to generalize to other types of KGs for structure and semantic enhancements. *Therefore, establishing a unified representation learning method for multiple types of KGs is worth to investigate.*

To overcome the challenges mentioned above, we propose a **Unified Hierarchical Representation learning method (UniHR)**, which includes a **Hierarchical Data Representation (HiDR)** module and a **Hierarchical Structure Learning (HiSL)** module as the graph encoder. HiDR module standardizes hyper-relational facts, nested factual facts, and temporal facts into the form of triples without loss of information. Furthermore, HiSL module captures local semantic information during intra-fact message passing and then utilizes inter-fact message passing to enrich the global structure information to obtain better node embeddings based on HiDR form. Finally, the updated embeddings are fed into decoders for link prediction. Empirical results demonstrate the UniHR achieves state-of-the-art performance on most metrics and highlight the strong potential of unified representations. Our contributions can be summarized as follows.

1. We emphasize the value of investigating unified KG representation method, including unified symbolic representation and unified representation learning method for different KGs.
2. To our knowledge, we propose the first unified KG representation learning framework UniHR, across different types of KGs, including a hierarchical data representation module and a hierarchical structure learning module.
3. We conduct link prediction experiments on 10 datasets across 5 types of KGs. Compared to methods designed for one kind of KG, UniHR achieves the best or competitive results, verifying strong generalization capability.

2 Preliminaries

In this section, we introduce the definition of four types of existing knowledge graphs (KGs): triple-

based KG, hyper-relational KG, nested factual KG and temporal KG, along with link prediction tasks on these types of KGs.

Link Prediction on Triple-based KG. A triple-based KG $\mathcal{G}_{KG} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}\}$ represents facts as triples, denoted as $\mathcal{F} = \{(h, r, t) \mid h, t \in \mathcal{V}, r \in \mathcal{R}\}$, where \mathcal{V} is the set of entities and \mathcal{R} is the set of relations. The link prediction on triple-based KGs involves answering a query $(h, r, ?)$ or $(?, r, t)$, where the missing element “?” is an entity in \mathcal{V} .

Link Prediction on Hyper-relational KG. A hyper-relational KG (HKG) $\mathcal{G}_{HKG} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}\}$ consists of hyper-relational facts, abbreviated as H-Facts, denoted as $\mathcal{F} = \{((h, r, t), \{(k_i: v_i)\}_{i=1}^m) \mid h, t, v_i \in \mathcal{V}, r, k_i \in \mathcal{R}\}$. Typically, we refer to (h, r, t) as the main triple in the H-Fact and $\{(k_i: v_i)\}_{i=1}^m$ as m auxiliary key-value pairs. Similar to link prediction on triple-based KGs, the link prediction on HKGs aims to predict entities in the main triple or the key-value pairs. Symbolically, the aim is to predict the missing element, denoted as “?” for queries $((h, r, t), (k_1: v_1), \dots, (k_i: ?))$, $((?, r, t), \{(k_i: v_i)\}_{i=1}^m)$ or $((h, r, ?), \{(k_i: v_i)\}_{i=1}^m)$.

Link Prediction on Nested Factual KG. A nested factual KG (NKG) can be represented as $\mathcal{G}_{NKG} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}, \hat{\mathcal{R}}, \hat{\mathcal{F}}\}$, which is composed of two levels of facts, called atomic facts and nested facts. $\mathcal{F} = \{(h, r, t) \mid h, t \in \mathcal{V}, r \in \mathcal{R}\}$ is the set of atomic facts, where \mathcal{V} is a set of atomic entities and \mathcal{R} is a set of atomic relations. $\hat{\mathcal{F}} = \{(\mathcal{F}_i, \hat{r}, \mathcal{F}_j) \mid \mathcal{F}_i, \mathcal{F}_j \in \mathcal{F}, \hat{r} \in \hat{\mathcal{R}}\}$ is the set of nested facts, where $\hat{\mathcal{R}}$ is the set of nested relations. The link prediction on the NKGs is performed on the atomic facts or nested facts. We refer to the link prediction on atomic facts as *Base Link Prediction*, and the link prediction on nested facts as *Triple Prediction*. For base link prediction, given a query $(h, r, ?)$ or $(?, r, t)$, the aim is to predict the missing atomic entity “?” from \mathcal{V} . For triple prediction, given a query $(?, \hat{r}, \mathcal{F}_j)$ or $(\mathcal{F}_i, \hat{r}, ?)$, the aim is to predict the missing atomic fact “?” from \mathcal{F} .

Link Prediction on Temporal KG. A temporal KG (TKG) $\mathcal{G}_{TKG} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}, \mathcal{T}\}$ is composed of quadruple-based facts, which can be represented as $\mathcal{F} = \{(h, r, t, [\tau_b, \tau_e]) \mid h, t \in \mathcal{V}, r \in \mathcal{R}, \tau_b, \tau_e \in \mathcal{T}\}$, where τ_b is the begin time, τ_e is the end time, \mathcal{V} is the set of entities, \mathcal{R} is the set of relations and \mathcal{T} is the set of timestamps. The link prediction on TKGs aims to predict missing entities “?”

in \mathcal{V} for two types of queries $(?, r, t, [\tau_b, \tau_e])$ or $(h, r, ?, [\tau_b, \tau_e])$.

3 Related Works

Link Prediction on Hyper-relational Knowledge Graph. Earlier HKG representation learning methods e.g., m-TransH (Wen et al., 2016), RAE (Zhang et al., 2018) have generalized the triple-based approach to HKG and loosely represent the combinations of key-value pairs. Galkin et al. first customize StarE (Galkin et al., 2020) based on CompGCN (Vashishth et al., 2019) for H-Facts to capture global structure information between H-Facts in the message passing stage, and achieves impressive results, demonstrating that the structure information of the graph in HKGs is also important. GRAN (Guan et al., 2021) introduces edge-aware bias into the vanilla transformer attention (Vaswani et al., 2017), while HyNT (Chung et al., 2023) designs a qualifier encoder for HKG. They both focus on intra-fact semantic dependencies but ignore the global structure information. Due to the existence of its particular key-value pairs on H-Facts, there are many limitations in capturing both global structure and local semantics.

Link Prediction on Nested Factual Knowledge Graph. Chung et al. (Chung and Whang, 2023) first introduced *nested facts* to imply relationships between facts. They also propose BiVE which bridges semantics between atomic facts and fact nodes in the encoding phase via a simple MLP and scores both atomic facts and nested facts using the quaternion-based KGE scoring functions like QuatE (Zhang et al., 2019) or BiQUE (Guo and Kok, 2021). Based on BiVE, NestE (Xiong et al., 2024) represents the fact nodes as a 1×3 embedding matrix and the nested relations as a 3×3 matrix to avoid information loss, embedding them into hyperplanes with different dimensions. These methods only capture the semantic information between atomic facts and nested facts while ignoring global structural information. Meanwhile, due to the complexity of this representation, common triple-based GNNs have difficulty in message passing between atomic fact and nested fact.

Link Prediction on Temporal Knowledge Graph. Recent studies in temporal knowledge graph representation learning have focused on enhancing performance by designing special time-aware scoring functions. Models such as TTransE (Leblay and

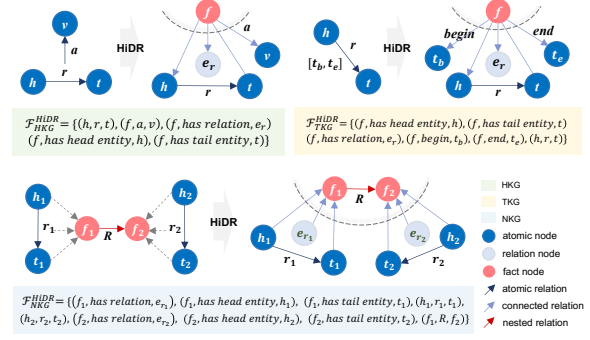


Figure 2: Diverse facts are translated into HiDR form.

Chekol, 2018), HyTE (Dasgupta et al., 2018), TeRo (Xu et al., 2020), and TGeomE+ (Xu et al., 2023a) incorporate temporal-aware module into the KGE score function in various ways. However, existing models seldom regard timestamps as node and directly utilize GNNs to perceive time information for enhancing entity and relation embeddings.

4 Methodology

In this section, we introduce our method, a **Unified Hierarchical Representation learning framework (UniHR)**, which includes a **Hierarchical Data Representation (HiDR)** module and a **Hierarchical Structure Learning (HiSL)** module. Our workflow can be divided into the following three steps: **1)** Given a KG \mathcal{G} of any type, we represent it into \mathcal{G}^{HiDR} under the HiDR form. **2)** The \mathcal{G}^{HiDR} will be encoded by HiSL module to enhance the semantic information within individual facts and structural information between facts on the whole graph. **3)** In the phase of decoding, the updated embeddings of nodes and edges are fed into transformer decoders to obtain the plausibility score of facts.

4.1 Hierarchical Data Representation

To overcome the differences in the representation of multiple types KGs, we introduce a **Hierarchical Data Representation** module, abbreviated as **HiDR**. Different from labelled RDF representation (Ali et al., 2022), we constrain “triple” to be considered as the basic units of HiDR form, then HiDR could continuous benefit from the model developments of triple-based KGs, which is the most active area about link prediction over KGs.

As shown in Fig. 2, in order to ensure comprehensive representation of facts, we introduce three hierarchical types of nodes and three connected relations in HiDR. Firstly, we denote original entities within three types of KGs as *atomic*

nodes and complement fact nodes for HKGs and TKGs lacking a designated fact node. To facilitate the interaction between fact nodes and relations explicitly, we incorporate *relation nodes* into the graph, represented as e_r for each r . These relation nodes are derived from transforming the relation edges in the original KG. It is important to facilitate direct access of fact nodes to the relevant atomic nodes during message passing process. To achieve this, we introduce three *connected relations*: has relation, has head entity and has tail entity, which establish directly connections between atomic nodes and fact nodes. Ultimately, we denote the (main) triple (h, r, t) in original fact as three *connected facts*: $(f, \text{has relation}, e_r)$, $(f, \text{has head entity}, h)$, $(f, \text{has tail entity}, t)$, and an *atomic fact* (h, r, t) , where f is fact node. Formally, the definition of HiDR form is as follows:

Definition 1. Hierarchical Data Representation:
A KG represented as the HiDR form is denoted as $\mathcal{G}^{HiDR} = \{\mathcal{V}^{HiDR}, \mathcal{R}^{HiDR}, \mathcal{F}^{HiDR}\}$, where $\mathcal{V}^{HiDR} = \mathcal{V}_a \cup \mathcal{V}_r \cup \mathcal{V}_f$ is a joint set of atomic node set (\mathcal{V}_a), relation node set (\mathcal{V}_r), fact node set (\mathcal{V}_f). $\mathcal{R}^{HiDR} = \mathcal{R}_a \cup \mathcal{R}_n \cup \mathcal{R}_c$ is a joint set of atomic relation set (\mathcal{R}_a), nested relation set (\mathcal{R}_n), connected relation set $\mathcal{R}_c = \{\text{has relation}, \text{has head entity}, \text{has tail entity}\}$. The fact set $\mathcal{F}^{HiDR} = \mathcal{F}_a \cup \mathcal{F}_c \cup \mathcal{F}_n$ is jointly composed of three types of triple-based facts: atomic facts (\mathcal{F}_a), connected facts (\mathcal{F}_c) and nested facts (\mathcal{F}_n), where $\mathcal{F}_a = \{(v_1, r, v_2) | v_1, v_2 \in \mathcal{V}_a, r \in \mathcal{R}_a\}$, $\mathcal{F}_c = \{(v_1, r, v_2) | v_1 \in \mathcal{V}_f, r \in \mathcal{R}_c, v_2 \in \mathcal{V}_a\}$, $\mathcal{F}_n = \{(v_1, r, v_2) | v_1, v_2 \in \mathcal{V}_f, r \in \mathcal{R}_n\}$.

Next, we introduce how to transform different types of KGs into HiDR form.

For hyper-relational knowledge graphs, we regard key-value pairs as complementary information for facts. Thus, we translate H-Facts $\mathcal{F}_{HKG} = \{((h, r, t), \{(k_i: v_i)\}_{i=1}^m)\}$ into the HiDR form that $\mathcal{G}_{HKG}^{HiDR} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}_{HKG}^{HiDR}\}$ following the definition, where $\mathcal{F}_c = \{(f, \text{has relation}, e_r), (f, \text{has head entity}, h), (f, \text{has tail entity}, t), (f, k_1, v_1), \dots, (f, k_m, v_m)\}$, $\mathcal{F}_a = \{(h, r, t) | ((h, r, t), \{(k_i: v_i)\}_{i=1}^m) \in \mathcal{F}_{HKG}\}$ and $\mathcal{F}_n = \emptyset$.

For nested factual knowledge graphs, HiDR can naturally represent hierarchical facts, so we translate the atomic facts $\mathcal{F}_{NKG} = \{(h_i, r_i, t_i)\}$ and the nested facts $\hat{\mathcal{F}}_{NKG} = \{((h_1, r_1, t_1), R, (h_2, r_2, t_2)) | (h_i, r_i, t_i) \in \mathcal{F}_{NKG}\}$ into the form of HiDR that $\mathcal{G}_{NKG}^{HiDR} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}_{NKG}^{HiDR}\}$ following the

definition, where $\mathcal{F}_a = \{(h_i, r_i, t_i) | (h_i, r_i, t_i) \in \mathcal{F}_{NKG}\}$, $\mathcal{F}_c = \{(f_i, \text{has head entity}, h_i), (f_i, \text{has tail entity}, t_i), (f_i, \text{has relation}, e_{r_i}) | f_i = (h_i, r_i, t_i) \in \mathcal{F}_{NKG}\}$ and $\mathcal{F}_n = \{(f_1, R, f_2) | f_i \in \mathcal{F}_{NKG}\}$.

For temporal knowledge graphs, we regard the TKG as a special HKG, and convert timestamps to auxiliary key-value pairs in HKGs by adding two special *atomic relations*: begin and end, regarding timestamps as special numerical atomic nodes. Thus, we firstly translate the temporal facts in TKGs $\mathcal{F}_{TKG} = \{(h, r, t, [\tau_b, \tau_e])\}$ into H-Facts form $\mathcal{F}_{TKG}^{HKG} = \{(h, r, t, \text{begin}:\tau_b, \text{end}:\tau_e)\}$. Then according to the previous transformation in HKG, it can be translated into the HiDR form that $\mathcal{G}_{TKG}^{HiDR} = \{\mathcal{V}, \mathcal{R}, \mathcal{F}_{TKG}^{HiDR}\}$ following the definition, where $\mathcal{F}_a = \{(h, r, t) | (h, r, t, \text{begin}:\tau_b, \text{end}:\tau_e) \in \mathcal{F}_{TKG}^{HKG}\}$, $\mathcal{F}_c = \{(f, \text{has relation}, e_r), (f, \text{has head entity}, h), (f, \text{has tail entity}, t), (f, \text{begin}, \tau_b), (f, \text{end}, \tau_e) | f = (h, r, t, \text{begin}:\tau_b, \text{end}:\tau_e) \in \mathcal{F}_{TKG}^{HKG}\}$ and $\mathcal{F}_n = \emptyset$.

In summary, we could convert all above KGs into the HiDR form, and preserve the semantics in the original KGs without loss of information.

4.2 Hierarchical Structure Learning

It's evident that HiDR form introduces many additional relation nodes and fact nodes. To avoid significantly increasing the model's training parameters while fully capturing the hierarchy of HiDR form, we design a Hierarchical Structure Learning module, abbreviated as **HiSL** shown in Fig. 3.

Representation Initialization. We first initialize the embedding matrices $\mathbf{H}_a \in \mathbb{R}^{|\mathcal{V}_a| \times d}$ and $\mathbf{E} \in \mathbb{R}^{|\mathcal{R}^{HiDR}| \times d}$ for atomic nodes and all relation edges. Then we also initialize the embedding of relation node $\mathbf{H}_r \in \mathbb{R}^{|\mathcal{V}_r| \times d}$, which can be transformed from the relation edge r with a projection matrix $\mathbf{W}_r \in \mathbb{R}^{d \times d}$: $\mathbf{H}_r = \mathbf{E}_a \cdot \mathbf{W}_r$, where $\mathbf{E}_a \subseteq \mathbf{E}$ is the atomic relation embeddings. Then we initialize the fact node embeddings \mathbf{H}_f to explicitly capture key information within facts by utilizing the embedding of (main) triple:

$$\mathbf{h}_f = f_m([\mathbf{h}_h; \mathbf{h}_r; \mathbf{h}_t]), \quad (1)$$

where $(h, r, t) \in \mathcal{F}_a$, $[\cdot; \cdot]$ is the concatenation operation, $\mathbf{h}_h, \mathbf{h}_t \subseteq \mathbf{H}_a$, $\mathbf{h}_r \subseteq \mathbf{H}_r$ denote (main) triple embedding and $f_m: \mathbb{R}^{3d} \rightarrow \mathbb{R}^d$ is a 1-layer MLP. Therefore, the initialization of relation nodes and fact nodes is sufficiently parameter-efficient.

For numerical atomic nodes, namely timestamps in temporal knowledge graphs, we encode the

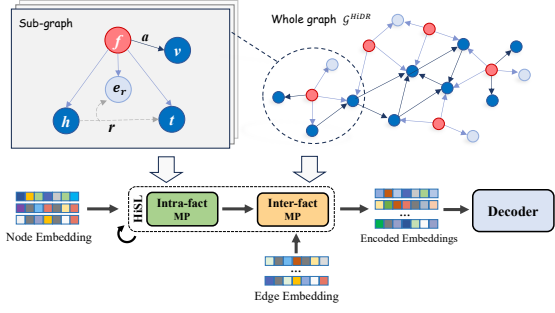


Figure 3: HiSL module for intra-fact and inter-fact MP.

timestamp τ into an embedding with Time2Vec (Kazemi et al., 2019):

$$\mathbf{h}_\tau = \omega_p \sin(f_p(\tau)) + f_{np}(\tau), \quad (2)$$

where $f_p: \mathbb{R}^1 \rightarrow \mathbb{R}^d$ is a 1-layer MLP as periodic function, $f_{np}: \mathbb{R}^1 \rightarrow \mathbb{R}^d$ is a 1-layer MLP as non-periodic function, and $\omega_p \in \mathbb{R}^1$ is a learnable parameter for scaling the periodic features.

Intra-fact Message Passing. In this stage, message passing is conducted for fact nodes. Given a fact node $f_k \in \mathcal{V}_f$, we construct its constituent elements, i.e., one-hop neighbors, as the node set $\mathcal{V}_k = \{v \in \mathcal{N}_{f_k} \mid v \in \mathcal{V}_a \cup \mathcal{V}_r\}$, where \mathcal{N}_{f_k} is the set of one-hop neighbors of fact node f_k . Then we retain the edges directly connected to fact node f_k , thereby constructing a subgraph $\mathcal{G}_k = \{\mathcal{V}_k, \mathcal{R}_k, \mathcal{F}_k\} \subseteq \mathcal{G}^{HiDR}$. For this subgraph, we employ the graph attention network (Brody et al., 2021) to aggregate local information, computing the attention score $\alpha_{i,j}$ between node $i \in \mathcal{V}_k$ and its neighbor j . The formula for calculating $\alpha_{i,j}$ in the l -th layer is as follows:

$$\alpha_{i,j}^l = \frac{\exp(\mathbf{W}^l(\sigma(\mathbf{W}_{in}^l \mathbf{h}_i^l + \mathbf{W}_{out}^l \mathbf{h}_j^l)))}{\sum_{j' \in \mathcal{N}_i} \exp(\mathbf{W}^l(\sigma(\mathbf{W}_{in}^l \mathbf{h}_i^l + \mathbf{W}_{out}^l \mathbf{h}_{j'}^l)))}, \quad (3)$$

where $\mathbf{h}_i^l, \mathbf{h}_j^l \in \mathbb{R}^d$ represent the embeddings of node i and its neighbor j in l -th layer. And there are three learnable weight matrices $\mathbf{W}_{in}^l, \mathbf{W}_{out}^l \in \mathbb{R}^{d \times d}$ and $\mathbf{W}^l \in \mathbb{R}^d$. We choose LeakyReLU as activation function σ . Then, the updated node embeddings are obtained by aggregating the information of neighbors according to the attention scores:

$$\mathbf{h}_i^l = \mathbf{h}_i^l + \sum_{j \in \mathcal{N}_i} \alpha_{i,j}^l \cdot \mathbf{W}_{out}^l \mathbf{h}_j^l. \quad (4)$$

Inter-fact Message Passing. At this stage, message passing is conducted on the whole graph \mathcal{G}^{HiDR} . Similar to previous work (Vashishth et al., 2019), we use a non-parametric aggregation operator $\phi(\cdot): \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ to obtain messages

from neighbouring nodes and edges. We employ the circular-correlation operator, defined as:

$$\phi(\mathbf{h}_j, \mathbf{e}_r) = \mathbf{h}_j \star \mathbf{e}_r = \mathbf{F}^{-1}((\mathbf{F}\mathbf{h}_j) \odot \overline{(\mathbf{F}\mathbf{e}_r)}) \quad (5)$$

where \mathbf{F} and \mathbf{F}^{-1} denote the discrete fourier transform (DFT) matrix and its inverse matrix, the \odot is the element-wise (Hadamard) product. Furthermore, in order to fully capture the heterogeneity of the graph, we classify edges along two dimensions: $\lambda(r) \in \{\text{forward}, \text{reverse}\}$ and $\tau(r) \in \{\text{connected relation}, \text{atomic relation}, \text{nested relation}\}$ and adopt two relation-type specific learnable parameters $\mathbf{W}_{\lambda(r)} \in \mathbb{R}^{d \times d}$ and $\omega_{\tau(r)} \in \mathbb{R}^1$ for more fine-grained aggregation:

$$\mathbf{h}_i^{l+1} = \sum_{(r,j) \in \mathcal{N}(i)} \sigma(\omega_{\tau(r)}^l) \mathbf{W}_{\lambda(r)}^l \phi(\mathbf{h}_j^l, \mathbf{e}_r^l) + \mathbf{W}_{self}^l \mathbf{h}_i^l \quad (6)$$

$$\mathbf{e}_r^{l+1} = \mathbf{W}_{rel}^l \mathbf{e}_r^l \quad (7)$$

where $\mathbf{W}_{self}^l, \mathbf{W}_{rel}^l \in \mathbb{R}^{d \times d}$, σ is a sigmoid activation function and $\mathcal{N}(i)$ is a set of immediate neighbors of i for its outgoing edges r . We utilize $\phi(\cdot)$ to combine the information from edge r and node j , and then passes it to node i for update.

Through Intra-fact and Inter-fact two-stage message passing, nodes can fully capture both local semantic and global structural information. Moreover, its number of training parameters does not increase with the scale of the graph, thereby effectively adapting to the HiDR form.

4.3 Link Prediction Decoder

Since the query varies across different settings, we use the transformer (Vaswani et al., 2017) as the decoder with the *mask* pattern. Specifically, we convert the updated node and edge embeddings into a sequence of fact embeddings, mask the elements to be predicted in facts with the $[M]$ token as the input to the transformer. Finally, we obtain the embedding of output $[M]$ in the last layer to measure the plausibility of the fact, denoted as \mathbf{h}_{pre} , and calculate the probability distribution of candidates, followed by training it using the cross-entropy loss:

$$P = \text{Softmax}(f(\mathbf{h}_{pre})[\mathbf{E}; \mathbf{H}]^T), \quad (8)$$

$$\mathcal{L} = \sum_{t=0}^{|\mathcal{R}|+|\mathcal{V}|} y_t \log P_t \quad (9)$$

where $P \in \mathbb{R}^{|\mathcal{R}|+|\mathcal{V}|}$ represents the confidence scores of all candidates, $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a 1-layer MLP, and $[\mathbf{E}; \mathbf{H}] \in \mathbb{R}^{(|\mathcal{R}|+|\mathcal{V}|) \times d}$ is the embedding matrix of all candidate edges or nodes. The P_t and y_t are probability and ground truth of the t -th candidate. The final loss function \mathcal{L} includes both node loss and edge loss during the predictions.

Model	WikiPeople						WD50K					
	subject/object			all entities			subject/object			all entities		
	MRR	H@1	H@10	MRR	H@1	H@10	MRR	H@1	H@10	MRR	H@1	H@10
NaLP	0.356	0.271	0.499	0.360	0.275	0.503	0.230	0.170	0.347	0.251	0.187	0.375
tNaLP	0.358	0.288	0.486	0.361	0.290	0.490	0.221	0.163	0.331	0.243	0.182	0.360
RAM	0.459	0.384	0.584	0.461	0.386	0.585	0.276	0.210	0.399	0.296	0.232	0.416
HINGE	0.393	0.309	0.547	0.395	0.311	0.549	0.264	0.187	0.410	0.277	0.200	0.424
NeuInfer	0.357	0.247	0.533	0.357	0.248	0.532	0.220	0.154	0.347	0.225	0.158	0.355
StarE	0.458	0.364	<u>0.611</u>	-	-	-	0.309	0.234	0.452	-	-	-
HyTransformer	0.460	0.382	0.594	-	-	-	0.304	0.231	0.443	-	-	-
GRAN	0.462	0.366	0.610	0.465	0.371	<u>0.613</u>	0.330	0.255	0.472	<u>0.361</u>	0.286	<u>0.501</u>
HyNT	<u>0.482</u>	<u>0.415</u>	0.602	<u>0.481</u>	<u>0.414</u>	0.603	<u>0.333</u>	<u>0.259</u>	<u>0.474</u>	<u>0.360</u>	<u>0.287</u>	0.500
UniHR	0.491	0.417	0.618	0.493	0.420	0.621	0.348	0.278	0.482	0.382	0.313	0.513

Table 1: Results of link prediction on HKG datasets. The best results are **bold**, while the second are underlined.

5 Experiment

5.1 Experiment Settings

Datasets. For link prediction on HKGs, we select three benchmark datasets: WikiPeople (Guan et al., 2021), WD50K (Galkin et al., 2020) and JF17K (Wen et al., 2016). As for the NKGs, we choose FBH, FBHE and DBHE constructed by (Chung and Whang, 2023). Lastly, we use wikidata12k (Dasgupta et al., 2018), a subset of wikidata (Vrandečić and Krötzsch, 2014) for link prediction on TKGs. The statistics of datasets are given in Appendix D.

Evaluation Metric. We conduct link prediction across multiple settings, evaluating performance based on the rank of predicted facts. We use the MR (Mean Rank), MRR (Mean Reciprocal Rank) and Hits@K (K=1,3,10) as our evaluation metrics. We abbreviate ‘Hits@K’ as ‘H@K’ and employ filtering settings (Bordes et al., 2013) during the evaluation to eliminate existing facts in the dataset.

Baselines. For link prediction on HKG, we compare our UniHR against NaLP (Guan et al., 2021), tNaLP (Guan et al., 2021), RAM (Liu et al., 2021), HINGE (Rosso et al., 2020), NeuInfer (Guan et al., 2020), StarE (Galkin et al., 2020), HyTransformer (Yu and Yang, 2021), GRAN (Wang et al., 2021) and HyNT (Chung et al., 2023). For link prediction on NKG, QuatE (Zhang et al., 2019), BiQUE (Guo and Kok, 2021), Neural-LP (Yang et al., 2017), DRUM (Sadeghian et al., 2019), AnyBURL (Meilicke et al., 2019), BiVE (Chung and Whang, 2023) and NestE (Xiong et al., 2024) are chosen as baselines. BiVE and NestE are especially designed for NKG. We compare against following TKG link prediction methods: ComplEx-N3 (Lacroix et al., 2018), HyTE (Dasgupta et al., 2018), TADistMult (Garcia-Duran et al., 2018), ATiSE (Xu et al., 2019), TeRo (Xu et al., 2020), TGeomE+

(Xu et al., 2023a), HGE (Pan et al., 2024).

Implementation details. All experiments are conducted on a single Nvidia 80G A800 GPU and implemented with PyTorch. For base link prediction on NKGs, we also use augmented triples from (Chung and Whang, 2023) for training to ensure fairness. For triple prediction, due to the small size of training set, we conduct training based on fixed embeddings of entities obtained from the base link prediction and set $\omega_{nested\ relation}=0$ to prevent overfitting. We employ AdamW (Kingma and Ba, 2015) optimizer. Hyperparameter details can be found in Appendix E.

5.2 Main Results

Link prediction on HKG. We compare our method with previous methods on the WD50K and WikiPeople datasets shown in Table 1. Among these methods, it can be seen that our proposed UniHR achieves state-of-the-art results, which means our method effectively captures global structural information. Compared to the GNN-based method StarE, we achieve improvements of 3.9 points (12.6%) in MRR, 4.4 points (18.8%) in Hits@1 and 3.0 points (6.6%) in Hits@10 on WD50K. This indicates that the performance of StarE’s customized graph neural network is limited by its inability to flexibly capture key-value pair information. Moreover, since the embeddings for newly added fact nodes and relation nodes are computed from atomic facts, so our training parameters do not significantly increase.

Link Prediction on NKG. Our experiments on the NKGs consist of two tasks: base link prediction and triple prediction. From the results in Table 2, we can see that our proposed UniHR obtains competitive results as the first method to capture global structural information of NKGs. For base

Model	FBH		DBHE		FBH			FBHE			DBHE		
	MRR	H@10	MRR	H@10	MR	MRR	H@10	MR	MRR	H@10	MR	MRR	H@10
Base link prediction					Triple prediction								
QuatE	0.354	0.581	0.264	0.440	145603.8	0.103	0.114	94684.4	0.101	0.209	26485.0	0.157	0.179
BiQUE	0.356	0.583	0.274	0.446	81687.5	0.104	0.115	61015.2	0.135	0.205	19079.4	0.163	0.185
Neural-LP	0.315	0.486	0.233	0.357	115016.6	0.070	0.073	90000.4	0.238	0.274	21130.5	0.170	0.209
DRUM	0.317	0.490	0.237	0.359	115016.6	0.069	0.073	90000.3	0.261	0.274	21130.5	0.166	0.209
AnyBURL	0.310	0.526	0.220	0.364	108079.6	0.096	0.108	83136.8	0.191	0.252	20530.8	0.177	0.214
BiVE	0.370	0.607	0.274	0.422	6.20	0.855	0.941	8.35	0.711	0.866	3.63	0.687	0.958
NestE	0.371	0.608	0.289	0.443	3.34	0.922	0.982	3.05	0.851	0.962	2.07	0.862	0.984
UniHR	0.401	0.619	0.296	0.448	2.46	0.946	0.993	5.20	0.793	0.890	1.90	0.862	0.987

Table 2: Results of base link prediction (left) and triple prediction (right). All baselines’ results are taken from (Xiong et al., 2024). For BiVE (Chung and Whang, 2023) and NestE (Xiong et al., 2024), we pick their best variants.

Model	wikidata12k			
	MRR	H@1	H@3	H@10
ComplEx-N3	0.248	0.143	-	0.489
HyTE	0.253	0.147	-	0.483
TA-DistMult	0.230	0.130	-	0.461
TeRo	0.299	0.198	0.329	0.507
ATiSE	0.252	0.148	0.288	0.462
TGeomE+	0.333	0.232	0.361	0.546
HGE	0.290	0.176	0.323	0.514
UniHR	0.333	0.240	0.367	0.527

Table 3: Results of link prediction on wikidata12k.

link prediction task on triple-based KGs, UniHR achieves considerable improvements. Of particular note, the MRR of FBHE increases by 8.1%.

For triple prediction, we perform best on FBH and DBHE datasets, especially obtaining an improvement of 2.4 points in MRR on FBH, and achieve the second-best performance on FBHE, which suggests that structural information is also valuable for NKG and UniHR can effectively capture the heterogeneity of NKG to enhance node embeddings. In particular, as a unified method, we do not use the customized decoder for triples, while previous state-of-the-art methods do. We will further illustrate the effectiveness of UniHR equipped with other decoders in Appendix C.

Link Prediction on TKG. As shown in Table 3, we achieve competitive results on the wikidata12k, even surpassing TGeomE+ by 3.4% on Hits@1 and 1.7% on Hits@3. However, existing temporal knowledge graph embedding methods, such as TGeomE+, often employ time-aware decoders, which are challenging to generalize to other types of KGs. In contrast, our approach efficiently encodes timestamps as atomic nodes only during initialization and learns temporal information through message passing on graph structure, demonstrating that graph structure information is also beneficial for temporal knowledge graphs, highlighting the effectiveness of our UniHR.

Variant	FBH			DBHE		
	MR	MRR	H@10	MR	MRR	H@10
w/o initial \mathbf{h}_f	5.22	0.909	0.980	2.56	0.794	0.978
w/o \mathbf{W}_r	3.06	0.944	0.992	2.30	0.850	0.976
w/o intra-fact MP	3.97	0.897	0.972	2.02	0.842	0.983
w/o $\omega_{\tau(r)}$	2.70	0.934	0.992	2.69	0.810	0.973
w/o $\mathbf{W}_{\lambda(r)}$	2.50	0.941	0.992	2.37	0.810	0.978
w/o inter-fact MP	2.61	0.913	0.991	2.11	0.827	0.986
UniHR	2.46	0.946	0.993	1.90	0.862	0.987

Table 4: Ablation studies on triple prediction task.

5.3 Ablation Study on HiSL

We conduct ablation experiments on triple prediction, the most relevant task to fact nodes. As shown in Table 4, all variants with certain modules or parameters removed exhibit a decrease in performance. We can conclude that intra-fact and inter-fact message passing modules both play crucial roles, allowing UniHR to more fully represent the current fact node with enhanced structural information. We also change the initialization of fact node embeddings to learnable embeddings (w/o initial \mathbf{h}_f), and results indicate that initializing fact node representations is essential. It highlights key information in the facts, mitigating the noise introduced by excessive neighbors.

5.4 Potential of Unified Representation

Generalize to Hyper-relational TKGs. Due to the unity of representation, UniHR can easily generalize to compositional types of KGs, like hyper-relational TKGs (Ding et al., 2024) that integrate HKGs and TKGs. From results in Table 6, UniHR offers a significant performance improvement in link prediction tasks on hyper-relational TKGs compared to TKG and HKG models, and even achieves competitive performance with the specialized hyper-relational TKG model HypeTKG.

Joint Learning on Different Tasks of KGs. For link prediction on NKGs, the two subtasks, namely

Model	WikiPeople ⁻										wikidata12k ⁻				
	subject/object					all entities					subject/object				
	MR	MRR	H@1	H@3	H@10	MR	MRR	H@1	H@3	H@10	MR	MRR	H@1	H@3	H@10
UniHR	835.8	0.486	0.412	0.528	0.617	829.0	0.488	0.414	0.531	0.620	818.7	0.314	0.220	0.345	0.509
UniHR _{Joint}	692.7	0.488	0.409	0.533	0.629	686.5	0.490	0.414	0.536	0.632	489.5	0.315	0.222	0.346	0.498

Table 5: Results of separate training and joint training on the HKG and TKG dataset, where identical entities and relations share the same embeddings. WikiPeople⁻ and wikidata12k⁻ represent the filtered test sets.

Model	WiKi-hy				YAGO-hy			
	MRR	H@1	H@3	H@10	MRR	H@1	H@3	H@10
TeRo	0.572	0.473	0.640	0.727	0.760	0.720	0.782	0.822
HGE	0.602	0.507	0.666	0.765	0.790	0.760	0.814	0.837
StarE	0.565	0.491	0.599	0.703	0.765	0.737	0.776	0.820
GRAN	0.661	0.610	0.679	0.750	0.808	0.789	0.817	0.842
HyNT	0.537	0.444	0.587	0.723	0.763	0.724	0.787	0.836
HypeTKG	0.687	0.633	0.710	0.789	0.832	0.817	0.838	0.857
UniHR	<u>0.681</u>	<u>0.618</u>	<u>0.717</u>	<u>0.774</u>	<u>0.823</u>	<u>0.801</u>	<u>0.839</u>	<u>0.859</u>

Table 6: Results on hyper-relational TKG datasets.

base link prediction and triple prediction, share the same graph during the message-passing phase under our representation form. Therefore, we attempt joint training on the two tasks using the NKG dataset, as shown in Table 7. We observe that the results of joint training are generally superior to those of separate training, indicating that nested facts and atomic facts can mutually enhance the capture of semantics .

Joint Learning on Different Types of KGs. We suppose that unified representation makes it possible to develop pre-trained models that integrate multiple types of KGs. To explore its potential, we conduct joint learning on different types of KGs. Therefore, we construct a hybrid dataset called **wikimix** which includes two subsets of Wikidata (Vrandečić and Krötzsch, 2014), namely HKG dataset WikiPeople and TKG dataset wikidata12k, which encompass 3547 identical entities and 18 identical relations. Due to the different types of facts, there are no identical facts in these two subsets. To further prevent data leakage, we filter out 537 entries from the wikidata12k test set whose main triples appear in the H-Facts of WikiPeople train set, and 384 entries from the WikiPeople test set whose main triples appear in wikidata12k train set. Statics of wikimix are given in Appendix D.

From the results in Table 5, it is evident that joint learning outperforms separate learning across most metrics. Notably, there are improvements of 17.1% and 40.2% in MR metric on wikippeople⁻ and wikidata12k⁻ datasets, respectively. This indicates that more complex structural interactions and diverse types of training data are beneficial. Moreover, our UniHR demonstrates good scalability and

Model	FBH		DBHE		FBH		DBHE	
	MRR	H@10	MRR	H@10	MRR	H@10	MRR	H@10
	Base link prediction				Triple prediction			
UniHR	0.401	0.619	0.296	0.448	0.946	0.993	0.862	0.987
UniHR _{Joint}	0.402	0.622	0.298	0.452	0.948	0.993	0.857	0.989

Table 7: Results of separate and joint training on NKG.

Dataset	Model	Training Params	Training Times
WD50K	StarE	10.81M	48h+
	HyNT	29.61M	22.3h
	UniHR	10.55M	23.4h
wikidata12k	ATiSE	31.46M	5.8h
	HGE	3.44M	-
	UniHR	3.68M	3.6h
YAGO-hy	HypeTKG	10.83M	37.5h
	UniHR	4.37M	17.6h

Table 8: Number of training parameters and times.

effectiveness in integrating multiple types of KGs.

5.5 Efficiency Analysis

Compared to existing state-of-the-art methods, UniHR does not significantly increase the overhead in terms of model training parameters and runtime as shown in Figure 8. The embeddings of relation nodes and fact nodes are computed using a simple MLP based on atomic nodes and relations, which does not significantly increase the number of training parameters. For runtime, we do not use the whole graph during all message passing phases and adopt a dropout strategy to prevent overfitting, which can enhance training efficiency.

6 Conclusion

In this paper, we propose UniHR, a unified hierarchical knowledge graph representation learning framework consisting of a Hierarchical Data Representation (HiDR) module and a Hierarchical Structure Learning (HiSL) module. The HiDR form unifies the hyper-relational facts, nested facts and temporal facts into the form of triples, overcoming the limitations of customized encoders for different forms of facts. Moreover, HiSL captures local semantic information within facts and global structural information between facts. Our UniHR achieves the best or competitive performance across five types of KGs, and highlight the strong potential of unified representations.

Limitations

The limitations are summarized as follows:

Our UniHR In this paper, our UniHR framework focuses on link prediction tasks under transductive settings with a single modality. In the future, we will investigate how to generalize our HiDR form to more complex tasks such as inductive reasoning (Teru et al., 2020) and multi-modal scenarios, etc.

Joint Learning on Different Types of KGs Constrained by computational resources, our analysis of the potential for joint training across multiple types of knowledge graphs focus only on HKG and TKG. We believe the unification of knowledge graph representation learning methods is a developing trend that makes it possible to develop unified pre-trained models based on multiple types of KGs. In the future, we aim to explore joint training across more types of KG to demonstrate the advantages of integrating multiple types of KG data.

Ethics Statement

In this paper, we explore the unified knowledge graph link prediction problem, aiming to complete various types of knowledge graphs using a unified model with deep learning techniques. Our training and evaluation are based on publicly available and widely used datasets of different types of knowledge graphs. Therefore, we believe this does not violate any ethics.

References

Waqas Ali, Muhammad Saleem, Bin Yao, Aidan Hogan, and Axel-Cyrille Ngonga Ngomo. 2022. [A survey of RDF stores & SPARQL engines for querying knowledge graphs](#). *VLDB J.*, 31(3):1–26.

KM Annervaz, Somnath Basu Roy Chowdhury, and Ambedkar Dukkipati. 2018. Learning beyond datasets: Knowledge graph augmented neural networks for natural language processing. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 313–322.

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.

Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko.

2013. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26.

Shaked Brody, Uri Alon, and Eran Yahav. 2021. How attentive are graph attention networks? In *International Conference on Learning Representations*.

Chanyoung Chung, Jaeeun Lee, and Joyce Jiyoung Whang. 2023. [Representation learning on hyper-relational and numeric knowledge graphs with transformers](#). In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*, pages 310–322. ACM.

Chanyoung Chung and Joyce Jiyoung Whang. 2023. [Learning representations of bi-level knowledge graphs for reasoning beyond link prediction](#). In *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 4208–4216. AAAI Press.

Shib Sankar Dasgupta, Swayambhu Nath Ray, and Partha Talukdar. 2018. Hyte: Hyperplane-based temporally aware knowledge graph embedding. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 2001–2011.

Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. 2018. [Convolutional 2d knowledge graph embeddings](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 1811–1818. AAAI Press.

Zifeng Ding, Jingcheng Wu, Jingpei Wu, Yan Xia, Bo Xiong, and Volker Tresp. 2024. [Temporal fact reasoning over hyper-relational knowledge graphs](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 355–373. Association for Computational Linguistics.

Mikhail Galkin, Priyansh Trivedi, Gaurav Maheshwari, Ricardo Usbeck, and Jens Lehmann. 2020. Message passing for hyper-relational knowledge graphs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7346–7359.

Alberto Garcia-Duran, Sebastijan Dumančić, and Mathias Niepert. 2018. Learning sequence encoders for temporal knowledge graph completion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4816–4821.

738	Saiping Guan, Xiaolong Jin, Jiafeng Guo, Yuanzhuo Wang, and Xueqi Cheng. 2020. Neuinfer: Knowledge inference on n-ary facts. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 6141–6151.	794
739		795
740		796
741		797
742		798
		799
743	Saiping Guan, Xiaolong Jin, Jiafeng Guo, Yuanzhuo Wang, and Xueqi Cheng. 2021. Link prediction on n-ary relational data based on relatedness evaluation. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 35(1):672–685.	800
744		801
745		802
746		803
747		804
748	Jia Guo and Stanley Kok. 2021. Bique: Biquaternionic embeddings of knowledge graphs. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8338–8351.	805
749		
750		
751		
752	Magdalena Kaiser, Rishiraj Saha Roy, and Gerhard Weikum. 2021. Reinforcement learning from reformulations in conversational question answering over knowledge graphs. In <i>Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval</i> , pages 459–469.	806
753		807
754		808
755		809
756		810
757		811
		812
		813
		814
		815
758	Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus A. Brubaker. 2019. Time2vec: Learning a vector representation of time . <i>CoRR</i> , abs/1907.05321.	816
759		817
760		818
761		819
762		820
763	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	821
764		822
765		823
766		824
767		825
768	Timothée Lacroix, Nicolas Usunier, and Guillaume Obozinski. 2018. Canonical tensor decomposition for knowledge base completion. In <i>International Conference on Machine Learning</i> , pages 2863–2872. PMLR.	826
769		827
770		828
771		829
772		
773	Julien Leblay and Melisachew Wudage Chekol. 2018. Deriving validity time in knowledge graph. In <i>Companion proceedings of the the web conference 2018</i> , pages 1771–1776.	830
774		831
775		832
776		833
777	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N. Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, and Christian Bizer. 2015. Dbpedia - A large-scale, multilingual knowledge base extracted from wikipedia . <i>Semantic Web</i> , 6(2):167–195.	834
778		835
779		836
780		837
781		838
782		839
783	Ye Liu, Hui Li, Alberto García-Durán, Mathias Niepert, Daniel Oñoro-Rubio, and David S. Rosenblum. 2019. MMKG: multi-modal knowledge graphs . In <i>The Semantic Web - 16th International Conference, ESWC 2019, Portorož, Slovenia, June 2-6, 2019, Proceedings</i> , volume 11503 of <i>Lecture Notes in Computer Science</i> , pages 459–474. Springer.	840
784		
785		
786		
787		
788		
789		
790	Yu Liu, Quanming Yao, and Yong Li. 2021. Role-aware modeling for n-ary relational knowledge bases. In <i>Proceedings of the Web Conference 2021</i> , pages 2660–2671.	841
791		842
792		843
793		844
		845
		846
		847
		848
		849
	Christian Meilicke, Melisachew Wudage Chekol, Daniel Ruffinelli, and Heiner Stuckenschmidt. 2019. Any-time bottom-up rule learning for knowledge graph completion. In <i>Proceedings of the 28th International Joint Conference on Artificial Intelligence</i> , pages 3137–3143.	
	Deepak Nathani, Jatin Chauhan, Charu Sharma, and Manohar Kaul. 2019. Learning attention-based embeddings for relation prediction in knowledge graphs. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4710–4723.	
	Jiaxin Pan, Mojtaba Nayyeri, Yinan Li, and Steffen Staab. 2024. HGE: embedding temporal knowledge graphs in a product space of heterogeneous geometric subspaces . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada</i> , pages 8913–8920. AAAI Press.	
	Paolo Rosso, Dingqi Yang, and Philippe Cudré-Mauroux. 2020. Beyond triplets: hyper-relational knowledge graph embedding for link prediction. In <i>Proceedings of the web conference 2020</i> , pages 1885–1896.	
	Ali Sadeghian, Mohammadreza Armandpour, Patrick Ding, and Daisy Zhe Wang. 2019. Drum: End-to-end differentiable rule mining on knowledge graphs. <i>Advances in Neural Information Processing Systems</i> , 32.	
	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2018. Rotate: Knowledge graph embedding by relational rotation in complex space. In <i>International Conference on Learning Representations</i> .	
	Komal Teru, Etienne Denis, and Will Hamilton. 2020. Inductive relation prediction by subgraph reasoning. In <i>International Conference on Machine Learning</i> , pages 9448–9457. PMLR.	
	Kristina Toutanova and Danqi Chen. 2015. Observed versus latent features for knowledge base and text inference . In <i>Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality, CVSC 2015, Beijing, China, July 26-31, 2015</i> , pages 57–66. Association for Computational Linguistics.	
	Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. 2019. Composition-based multi-relational graph convolutional networks. In <i>International Conference on Learning Representations</i> .	
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	

850	Denny Vrandečić and Markus Krötzsch. 2014. Wiki-	907
851	data: a free collaborative knowledgebase. <i>Communi-</i>	908
852	cations of the ACM, 57(10):78–85.	909
853	Quan Wang, Haifeng Wang, Yajuan Lyu, and Yong Zhu.	910
854	2021. Link prediction on n-ary relational facts: A	911
855	graph-based approach. In <i>Findings of the Association</i>	912
856	<i>for Computational Linguistics: ACL-IJCNLP 2021</i> ,	913
857	pages 396–407.	914
858	Jianfeng Wen, Jianxin Li, Yongyi Mao, Shini Chen,	915
859	and Richong Zhang. 2016. On the representation	916
860	and embedding of knowledge bases beyond binary	917
861	relations. In <i>Proceedings of the Twenty-Fifth Inter-</i>	
862	<i>national Joint Conference on Artificial Intelligence</i> ,	
863	pages 1300–1307.	
864	Bo Xiong, Mojtaba Nayyeri, Daniel Daza, and Michael	
865	Cochez. 2023. Reasoning beyond triples: Recent	
866	advances in knowledge graph embeddings . In <i>Pro-</i>	
867	<i>ceedings of the 32nd ACM International Conference</i>	
868	<i>on Information and Knowledge Management, CIKM</i>	
869	<i>2023, Birmingham, United Kingdom, October 21-25,</i>	
870	<i>2023</i> , pages 5228–5231. ACM.	
871	Bo Xiong, Mojtaba Nayyeri, Linhao Luo, Zihao Wang,	
872	Shirui Pan, and Steffen Staab. 2024. Neste: Model-	
873	ing nested relational structures for knowledge graph	
874	reasoning . In <i>Thirty-Eighth AAAI Conference on</i>	
875	<i>Artificial Intelligence, AAAI 2024, Thirty-Sixth Con-</i>	
876	<i>ference on Innovative Applications of Artificial Intel-</i>	
877	<i>ligence, IAAI 2024, Fourteenth Symposium on Educa-</i>	
878	<i>tional Advances in Artificial Intelligence, EAAI 2014,</i>	
879	<i>February 20-27, 2024, Vancouver, Canada</i> , pages	
880	9205–9213. AAAI Press.	
881	Chengjin Xu, Mojtaba Nayyeri, Fouad Alkhoury,	
882	Hamed Shariat Yazdi, and Jens Lehmann. 2019. Tem-	
883	poral knowledge graph embedding model based on	
884	additive time series decomposition. <i>arXiv preprint</i>	
885	<i>arXiv:1911.07893</i> .	
886	Chengjin Xu, Mojtaba Nayyeri, Fouad Alkhoury,	
887	Hamed Shariat Yazdi, and Jens Lehmann. 2020. Tero:	
888	A time-aware knowledge graph embedding via tem-	
889	poral rotation. In <i>Proceedings of the 28th Inter-</i>	
890	<i>national Conference on Computational Linguistics</i> ,	
891	pages 1583–1593.	
892	Chengjin Xu, Mojtaba Nayyeri, Yung-Yu Chen, and	
893	Jens Lehmann. 2023a. Geometric algebra based	
894	embeddings for static and temporal knowledge	
895	graph completion . <i>IEEE Trans. Knowl. Data Eng.</i> ,	
896	35(5):4838–4851.	
897	Hongcai Xu, Junpeng Bao, and Wenbo Liu. 2023b.	
898	Double-branch multi-attention based graph neural	
899	network for knowledge graph completion. In <i>Pro-</i>	
900	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	
901	<i>tion for Computational Linguistics (Volume 1: Long</i>	
902	<i>Papers)</i> , pages 15257–15271.	
903	Fan Yang, Zhilin Yang, and William W Cohen. 2017.	
904	Differentiable learning of logical rules for knowledge	
905	base reasoning. <i>Advances in neural information pro-</i>	
906	<i>cessing systems</i> , 30.	
	Donghan Yu and Yiming Yang. 2021. Improving	907
	hyper-relational knowledge graph completion. <i>arXiv</i>	908
	<i>preprint arXiv:2104.08167</i> .	909
	Richong Zhang, Junpeng Li, Jiajie Mei, and Yongyi	910
	Mao. 2018. Scalable instance reconstruction in	911
	knowledge bases via relatedness affiliated embed-	912
	ding. In <i>Proceedings of the 2018 world wide web</i>	913
	<i>conference</i> , pages 1185–1194.	914
	Shuai Zhang, Yi Tay, Lina Yao, and Qi Liu. 2019.	915
	Quaternion knowledge graph embeddings. <i>Advances</i>	916
	<i>in neural information processing systems</i> , 32.	917
	A Related Works	918
	A.1 Link Prediction on Triple-based KGs	919
	Most existing techniques in KG representation	920
	learning are proposed for triple-based KGs. Among	921
	these techniques, knowledge graph embedding	922
	(KGE) models (Bordes et al., 2013; Sun et al.,	923
	2018) have received extensive attention due to their	924
	effectiveness and simplicity. The idea is to project	925
	entities and relations in the KG to low-dimensional	926
	vector spaces, utilizing KGE scoring functions to	927
	measure the plausibility of triples in the embedding	928
	space. Typical methods include TransE (Bordes	929
	et al., 2013), RotatE (Sun et al., 2018), and ConvE	930
	(Dettmers et al., 2018).	931
	Depending on the KGE model alone has lim-	932
	itation of capturing complex graph structures,	933
	whereas augmenting global structural information	934
	with a graph neural network (GNN) (Vashishth	935
	et al., 2019; Nathani et al., 2019; Xu et al., 2023b)	936
	proves to be an effective approach for enhance-	937
	ment. The paradigm of combining GNN as encoder	938
	with KGE scoring function as decoder helps to en-	939
	hance the performance of KGE scoring function.	940
	These GNN methods design elaborate message	941
	passing mechanisms to capture the global struc-	942
	tural features. Typically, CompGCN (Vashishth	943
	et al., 2019) aggregates the joint embedding of	944
	entities and relations in the neighborhood via a	945
	parameter-efficient way and MA-GNN (Xu et al.,	946
	2023b) learns global-local structural information	947
	based on multi-attention. These methods achieve	948
	impressive results on triple-based KGs but are hard	949
	to generalize to beyond-triple KGs.	950
	B Results on JF17K	951
	Table 9 shows the experimental results on JF17K.	952
	Due to the absence of a validation set in the JF17K	953
	dataset and the different ways of dividing the	954
	dataset across various baselines, we adopt the re-	955
	sults reported in the original paper. Consistent with	956

previous experiments on hyper-relational knowledge graphs, we also achieve state-of-the-art performance on JF17K among all baselines. In particular, we achieved 1.7 (2.9%) points in MRR and 1.5 (2.9%) points in Hits@1 compared to the method StarE which also utilises a graph neural network encoding and a simple transformer decoding, indicating that our hierarchical GNN HiSL could better capture the structure of hyper-relational facts.

Model	JF17K					
	subject/object			all entities		
	MRR	Hits@1	Hits@10	MRR	Hits@1	Hits@10
m-TransH	0.206	0.206	0.462	-	-	-
NaLP	0.221	0.165	0.331	0.366	0.290	0.516
HINGE	0.431	0.342	0.611	0.517	0.436	0.675
NeuInfer	0.449	0.361	0.624	0.473	0.397	0.618
StarE	0.574	0.496	0.725	-	-	-
HyTransformer	0.582	0.501	0.742	-	-	-
UniHR	0.591	0.511	0.745	0.621	0.545	0.768

Table 9: Link prediction on JF17K. All results of baselines are taken from the original paper. Best results are in bold.

C Decoder Analysis

To explore the effectiveness of our UniHR encoding further, we pair UniHR with different decoders and evaluated them on triple prediction task. In addition to the previously mentioned unified framework **UniHR + Transformer**, we also experiment on **UniHR + ConvE** with two scoring strategies. The ConvE (Dettmers et al., 2018) is the decoder customized for triples and its scoring function is $vec\left(\sigma\left(\left[\tilde{\mathbf{h}}_h; \tilde{\mathbf{e}}_r\right] * \psi\right)\right)$, where $\tilde{\mathbf{h}}_h$ and $\tilde{\mathbf{e}}_r$ represent reshaped 2D embeddings of head entity h and relation r , and $*$ is a convolution operator. The $vec(\cdot)$ and ψ are denoted as the vectorization function and a set of convolution kernels.

Due to our special representation, there exists two scoring methods for atomic triples, thus we present the base link prediction results separately for each scoring method. The s_f represents scoring triples $(f, has\ head\ entity, h)$ and $(f, has\ tail\ entity, t)$, and s_t represents scoring (h, r, t) . The performance of base link prediction is shown in Table 10. Notably, FBH and FBHE share identical atomic facts, resulting in the same performance. It can be observed that regardless of the scoring method employed, we both achieve competitive performance, especially with scoring (h, r, t) on FBH and scoring $(f, has\ head\ entity, h)$ $(f, has\ tail\ entity, t)$ on DBHE. We attribute the differences in performance

Model	FBHE/FBH		DBHE	
	MRR	Hits@10	MRR	Hits@10
QuatE	0.354	0.581	0.264	0.440
BiQUE	0.356	0.583	0.274	0.446
Neural-LP	0.315	0.486	0.233	0.357
DRUM	0.317	0.490	0.237	0.359
AnyBURL	0.310	0.526	0.220	0.364
BiVE	0.370	0.607	0.274	0.422
NestE	0.371	<u>0.608</u>	0.289	0.443
UniHR + ConvE s_h	<u>0.397</u>	0.622	0.289	0.443
UniHR + ConvE s_f	0.375	0.596	0.307	0.471
UniHR + Transformer	0.401	<u>0.619</u>	<u>0.296</u>	<u>0.448</u>

Table 10: Base link prediction on FBHE, FBH and DBHE. All baselines’ results are taken from (Xiong et al., 2024). The best results among all models are written bold, while the second are underlined. The s_f and s_h denote $(f, has\ head\ entity, h)$ $(f, has\ tail\ entity, t)$ and (h, r, t) two types of scoring method respectively. For BiVE (Chung and Whang, 2023) and NestE (Xiong et al., 2024), we pick their variants with best performance.

under different scoring methods to dataset characteristics. The DBHE dataset is relatively smaller, and scoring method s_f effectively alleviates overfitting problem. Conversely, for larger datasets FBH, scoring based on (h, r, t) minimizes information loss.

Table 11 shows the results of triple prediction on three benchmark datasets. Among all baselines, Quate, Bique, Neural-LP, Drum, and AnyBURL struggle to model the mapping relationship between atomic facts and nested facts. Furthermore, prior works (Chung and Whang, 2023) do not guarantee that all atomic facts in the nested fact test set are present in the training set as entities, which shifts the problem from a transductive setting to an inductive setting, leading to significant performance gaps between these baselines. On most metrics, our method outperforms BiVE and NestE which are specifically modeled for nested facts. Notably, NestE fully preserves the semantics of atomic facts. However, on the FBHE dataset, UniHR + ConvE achieves an improvement of 0.58 (6.4%) points in MRR and 0.24 (2.4%) points in Hits@10 compared to the state-of-the-art model NestE and the second-best performance after UniHR + Transformer on the FBH and DBHE datasets, demonstrating UniHR’s powerful graph structure encoding capabilities. We also carry out ablation experiments on UniHR + ConvE as shown in Table 11. Performance declines are observed after removing any part of the HiSL module, showing the significance of HiSL for hierarchical encoding.

Model	FBH			FBHE			DBHE		
	MR	MRR	Hits@10	MR	MRR	Hits@10	MR	MRR	Hits@10
QuatE (Zhang et al., 2019)	145603.8	0.103	0.114	94684.4	0.101	0.209	26485.0	0.157	0.179
BiQUE (Guo and Kok, 2021)	81687.5	0.104	0.115	61015.2	0.135	0.205	19079.4	0.163	0.185
Neural-LP (Yang et al., 2017)	115016.6	0.070	0.073	90000.4	0.238	0.274	21130.5	0.170	0.209
DRUM (Sadeghian et al., 2019)	115016.6	0.069	0.073	90000.3	0.261	0.274	21130.5	0.166	0.209
AnyBURL (Meilicke et al., 2019)	108079.6	0.096	0.108	83136.8	0.191	0.252	20530.8	0.177	0.214
BiVE (Chung and Whang, 2023)	6.20	0.855	0.941	8.35	0.711	0.866	3.63	0.687	0.958
NestE (Xiong et al., 2024)	3.34	0.922	0.982	3.05	0.851	0.962	2.07	0.862	0.984
UniHR + Transformer	2.46	0.946	0.993	5.20	0.793	0.890	1.90	0.862	0.987
UniHR + ConvE	3.00	0.900	0.983	6.27	0.909	0.986	2.06	0.876	0.978
UniHR + ConvE w/o h_f	4.39	0.887	0.979	10.10	0.865	0.970	2.76	0.798	0.961
UniHR + ConvE w/o intra-fact	6.54	0.859	0.959	18.10	0.871	0.968	5.82	0.665	0.900
UniHR + ConvE w/o inter-fact	12.56	0.864	0.961	20.56	0.864	0.966	10.75	0.764	0.951

Table 11: Triple prediction on FBHE, FBH and DBHE. All baselines’ results are taken from (Xiong et al., 2024). The best results among all models are written bold. For BiVE (Chung and Whang, 2023) and NestE (Xiong et al., 2024), we pick their variants with best performance.

Dataset	Fact	Entities	Relations	Train	Valid	Test	with Q(%)	Arity	Nested Fact	Nested Relation	with AF(%)	Period
<i>Hyper-relational Knowledge Graph</i>												
WikiPeople	369866	34825	178	294439	37715	37712	9482(2.6%)	2-7	-	-	-	-
WD50K	236507	47155	531	166435	23913	46159	32167(13.6%)	2-67	-	-	-	-
JF17K	100947	28645	501	76379	-	24568	46320(45.9%)	2-6	-	-	-	-
<i>Nested Factual Knowledge Graph</i>												
FBH	310116	14541	237	248094	31011	31011	-	-	27062	6	33157	-
FBHE	310116	14541	237	248094	31011	31011	-	-	34941	10	33719	-
DBHE	68296	12440	87	54636	6830	6830	-	-	6717	8	8206	-
<i>Temporal Knowledge Graph</i>												
wikidata12k	40621	12554	24	32497	4062	4062	-	-	-	-	-	19-2020
<i>Hyper-relational Temporal Knowledge Graph</i>												
Wiki-hy	139078	16634	147	111252	13900	13926	13335(9.59%)	2-8	-	-	-	1513-2020
YAGO-hy	73143	16167	54	51193	10973	10977	5107(6.98%)	2-5	-	-	-	0-187
<i>Multiple types of Knowledge Graph</i>												
wikimix	409566	43832	184	326936	41777	3525(TKG)/37328(HKG)	9098(2.2%)	2-7	-	-	-	19-2020

Table 12: The statistics of diverse knowledge graphs dataset, where “with Q(%)” and “Arity” column respectively denote the number of facts with auxiliary key-value pairs and the range of arity of hyper-relational facts, the “with AF(%)” column denotes the number of atomic facts in nested facts.

D Datasets Statistics

Table 12 shows the details of the three hyper-relational knowledge graph benchmark datasets: WikiPeople, WD50K, JF17K, three nested factual knowledge graph benchmark datasets: FBH, FBHE, DBHE, and the temporal knowledge graph benchmark dataset wikidata12k. Among them, WikiPeople is a dataset derived from Wikidata (Vrandečić and Krötzsch, 2014) concerning entities type “human”. WikiPeople filter out the elements which have at least 30 mentions as key-value pairs. WD50K is a high-quality dataset extracting from Wikidata statements and avoiding the potential data leakage which allows triple-based models to memorize main fact in the H-Facts of test set. The “with Q(%)” column in Table 12 denote the number of facts with auxiliary key-value pairs and the “Arity” column denote range of the number of entities in hyper-relational facts. The nested factual knowledge graph datasets FBH and FBHE (Chung and Whang, 2023) are constructed based

on FB15k237 (Toutanova and Chen, 2015) from Freebase (Bollacker et al., 2008) while DBHE is based on DB15K (Liu et al., 2019) from DBpedia (Lehmann et al., 2015). FBH contains nested facts that can be only inferred inside the atomic facts, while FBHE and DBHE contain externally-sourced nested relation crawling from Wikipedia articles, e.g., NextAlmaMater and SucceededBy. Temporal knowledge graph dataset wikidata12K is also a subset of Wikidata (Vrandečić and Krötzsch, 2014), which represents the time information $\tau \in \mathcal{T}$ as time intervals.

E Hyperparameter Settings

Here, we show the hyperparameter details for each link prediction task. To be specific, we tune the learning rate using the range $\{0.0001, 0.0005, 0.001\}$, the embedding dim using the range $\{50, 100, 200, 400\}$, the GNN layer using the range $\{1, 2, 3\}$ and dropout using the range $\{0.1, 0.2, 0.3, 0.4\}$. Additionally, we use

Hyperparameter	WikiPeople	WD50K	wikidata12k	FBHE _{base}	FBH _{base}	DBHE _{base}	FBHE _{triple}	FBH _{triple}	DBHE _{triple}
batch_size	2048	2048	2048	2048	2048	2048	2048	2048	2048
embedding dim	200	200	200	200	200	200	200	200	200
hidden dim	200	200	200	200	200	200	200	200	200
GNN_layer	2	2	2	2	2	2	2	2	2
GNN_intra-fact heads	4	4	4	4	4	4	4	4	4
GNN_intra-fact dropout	0.1	0.1	0.1	0.3	0.1	0.3	0.2	0.1	0.1
GNN_inter-fact activation	tanh	tanh	tanh	tanh	tanh	tanh	tanh	tanh	tanh
GNN_dropout	0.1	0.1	0.1	0.3	0.1	0.3	0.2	0.1	0.1
transformer layers	2	2	2	2	2	2	2	2	2
transformer heads	4	4	4	4	4	4	4	4	4
transformer activation	gelu	gelu	gelu	gelu	gelu	gelu	gelu	gelu	gelu
decoder dropout	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
soft label for entity	0.2	0.2	0.4	0.2	0.2	0.3	0.2	0.2	0.2
soft label for relation	0.1	0.1	0.3	0.2	0.2	0.3	0.2	0.2	0.2
weight_decay	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
learning rate	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4

Table 13: The major hyperparameters of our approach for all link prediction tasks.

smoothing label in the training phase from range
 $\{0.1, 0.2, 0.3\}$. The best hyperparameters obtained
from the experiments are presented in Table 13.