# Graphical Abstract
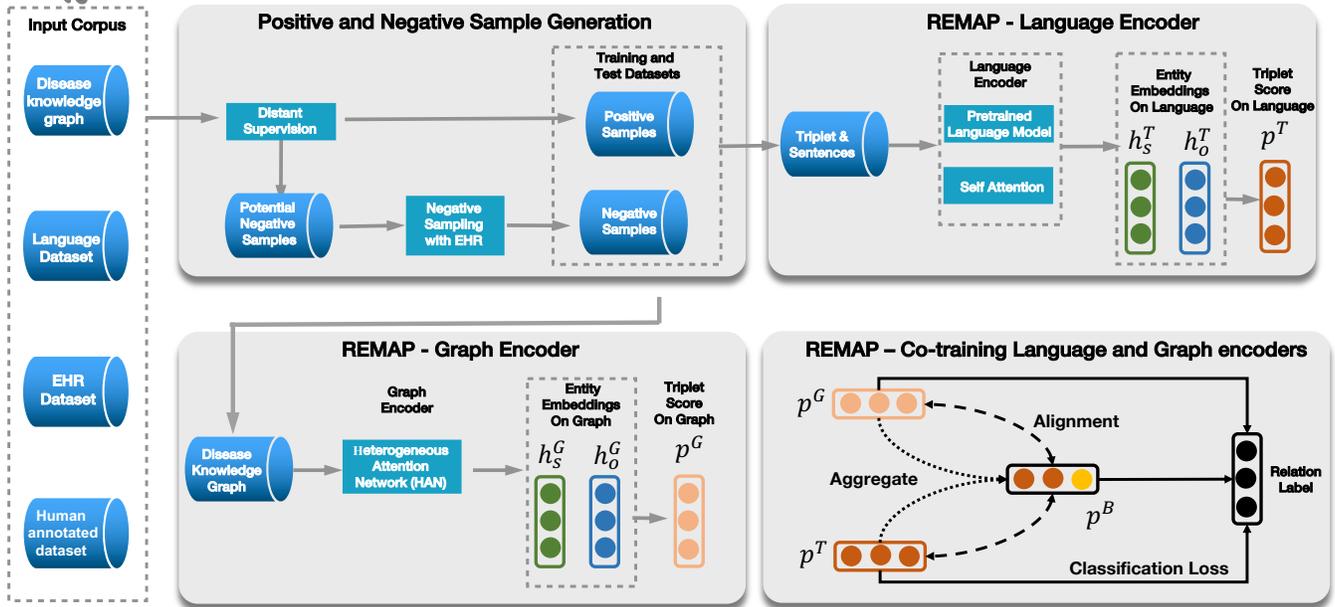
**Multimodal Learning on Graphs for Disease Relation Extraction**

Yucong Lin, Keming Lu, Sheng Yu, Tianxi Cai, Marinka Zitnik

## Multimodal Learning on Graphs for Disease Relation Extraction

# Highlights

**Multimodal Learning on Graphs for Disease Relation Extraction**

Yucong Lin, Keming Lu, Sheng Yu, Tianxi Cai, Marinka Zitnik

- We develop a flexible multimodal approach for extracting and classifying diverse kinds of disease-disease relationships (REMAP). REMAP fuses knowledge graph embeddings with deep language models and can flexibly accommodate missing data types.

- Evaluation against a clinical expert annotated dataset shows that REMAP achieves 88.6% micro-accuracy and 81.8% micro-F1 score, outperforming text-based methods by 10 and 17.2 percentage points, respectively.

- We release a high-quality test dataset of gold-standard annotations developed as a consensus of three clinical experts for evaluating disease-disease relation extraction, together with the open-source implementation of REMAP.

# Multimodal Learning on Graphs for Disease Relation Extraction

Yucong Lin[a,b,1], Keming Lu[c,1], Sheng Yu[d,e], Tianxi Cai[f,g], Marinka Zitnik[g,h,i,*]

[a]*Institute of Engineering Medicine, Beijing Institute of Technology, Beijing, China, Beijing, China*
[b]*Beijing Engineering Research Center of Mixed Reality and Advanced Display, School of Optics and Photonics,Beijing Institute of Technology, Beijing, China*
[c]*Viterbi School of Engineering, University of Southern California, Los Angeles, CA, 90007, USA*
[d]*Center for Statistical Science, Tsinghua University, Beijing, China*
[e]*Department of Industrial Engineering, Tsinghua University, Beijing, China*
[f]*Department of Biostatistics, Harvard T.H.Chan School of Public Health, Boston, MA, 02115, USA*
[g]*Department of Biomedical Informatics, Harvard Medical School, Boston, MA, 02115, USA*
[h]*Broad Institute of MIT and Harvard, Boston, MA, 02142, USA*
[i]*Harvard Data Science Initiative, Cambridge, MA, 02138, USA*

## Abstract

**Objective:** Disease knowledge graphs are a way to connect, organize, and access disparate information about diseases with numerous benefits for artificial intelligence (AI). To create knowledge graphs, it is necessary to extract knowledge from multimodal datasets in the form of relationships between disease concepts and normalize both concepts and relationship types.

**Methods:** We introduce REMAP, a multimodal approach for disease relation extraction and classification. The REMAP machine learning approach jointly embeds a partial, incomplete knowledge graph and a medical language dataset into a compact latent vector space, followed by aligning the multimodal embeddings for optimal disease relation extraction.

**Results:** We apply REMAP approach to a disease knowledge graph with 96,913 relations and a text dataset of 1.24 million sentences. On a dataset annotated by human experts, REMAP improves text-based disease relation extraction by 10.0% (accuracy) and 17.2% (F1-score) by fusing disease knowledge graphs with text information. Further, REMAP leverages text information to recommend new relationships in the knowledge graph, outperforming graph-based methods by 8.4% (accuracy) and 10.4% (F1-score).

**Conclusion:** REMAP is a multimodal approach for extracting and classifying disease relationships by fusing structured knowledge and text. REMAP provides a flexible neural architecture to easily find, access, and validate AI-driven relationships between disease concepts.

*Keywords:*
Disease relation extraction, Medical relation extraction, Multimodal learning, Knowledge graphs, Graph neural networks, Language neural models

## 1. Introduction

Disease knowledge graphs are a way to connect, organize, and access disparate data and information resources about diseases with numerous benefits for artificial intelligence (AI). Systematized

---

[*]Corresponding author: marinka@hms.harvard.com; 1-617-432-5138; Department of Biomedical Informatics, Harvard Medical School, 10 Shattuck Street, Boston, MA 02115, USA
[1]Equal contribution

knowledge can be injected into AI methods to imitate human experts' reasoning so that AI-driven hypotheses can be easily found, accessed, and validated. For example, disease knowledge graphs (KGs) power AI applications, such as identification of disease treatments [1] and electronic health record (EHR) retrieval [2]. However, creating high-quality knowledge graphs requires extracting relationships between diseases from disparate information sources, such as free text in the EHRs and semantic knowledge representation from literature.

Traditionally, KGs were constructed via manual efforts, requiring humans to input every fact [3]. In contrast, rule-based [4] and semi-automated [5, 6] methods, while scalable, can suffer from poor accuracy and low recall rates. As a result, an enticing alternative is to create KGs by extracting relationships from literature and building large-scale KGs that comprehensively cover a domain of interest. These methods leverage pre-trained language models [7, 8] and have advanced the analysis of biomedical knowledge graphs [9, 10, 11, 12, 13]. Another approach for populating KGs with relations uses knowledge graph embeddings (KGE), which directly predict new relations in partial, incomplete knowledge graphs. KGE methods learn how to represent every entity (i.e., node) and relation (i.e., edge) in a graph as a distinct point in a low-dimensional vector space (i.e., embedding) so that performing algebraic operations in this learned space reflects the topology of the graph [14]. Embeddings produced by KGE methods can be remarkably powerful for downstream AI applications [15, 16, 17, 18, 19, 20]. Widely used KGE methods include translation models [21, 22, 18, 23], bilinear models [24, 25, 26, 27], and graph neural networks (GNNs) [28, 29, 30, 31]. These methods leverage embeddings to predict new relations, thereby completing sparse knowledge areas and systematically growing an existing KG. However, extracting relations from a single data type may suffer from bias, noise, and incompleteness. For example, in language-based methods, the training dataset is collected using distant supervision [32], which creates noisy sentences that can mislead relation extraction. Further, graph-based methods can suffer from out-of-dictionary problems, which limit the ability to model relations involving entities previously not in the KG [21, 33].

Nevertheless, language-based and graph-based methods both have advantages. For example, language-based methods can reason over large datasets created using techniques such as distant supervision and contrastive learning, and graph-based methods can operate on noisy and incomplete knowledge graphs, providing robust predictions. An emerging strategy to advance relation extraction thus leverages multiple data types simultaneously [34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 22, 44, 45, 46, 47] with multimodal learning [48], outperforming rule-based [4] and semi-automated [5, 6] methods. However, existing approaches are limited in two ways, which we outline next.

First, KGs provide only positive samples (i.e., examples of true disease-disease relationships), while existing methods also require negative samples that, ideally, are disease pairs that resemble positive samples but are not valid relationships. Methods for positive-unlabeled [49] and contrastive [50, 51] learning can address this challenge by sampling random disease pairs from the dataset as negative proxy samples and ensuring a low false-positive rate. However, these methods may not generalize well in real-world applications because random negative samples are not necessarily realistic and fail to represent the boundary cases. To improve the quality of negative sampling in disease relation extraction, we introduce an EHR-based negative sampling strategy in this work. With the strategy, our approach generates negative samples using disease pairs that rarely appear together in EHRs, thus having realistic negative samples to enable the broad generalization of the approach.

Second, only graph or language information is available for a subset of diseases but not both modalities. For example, in Lin *et al.* [11], over 60% disease pairs in the KG had no corresponding text information; there were also cases with text but no graph information. Thus, multimodal approaches must be flexible, meaning they can make predictions when only one data type is available.

Unfortunately, a small number of existing multimodal approaches with such capability [52, 53, 54, 55] do not consider language and graphs. Further, some studies [56, 57] use adversarial learning to impute data from missing modalities, but the imputed values can introduce unwanted bias, leading to distribution shifts. To address this issue, in this work, we develop a multimodal de-coupled architecture where language and graph modules interact only through shared parameters and a cross-modal loss function. This approach ensures our model can take advantage of both language and graph inputs and identify disease relations using either single or multimodal inputs.

**Present work.** We introduce REMAP (Relation Extraction with Multimodal Alignment Penalty)[2], a multimodal approach for extracting and classifying disease-disease relations (Figure 1). REMAP is a flexible multimodal algorithm that jointly learns over text and graphs with a unique capability to make predictions even when a disease concept exists in only one data type. To this end, REMAP specifies graph-based and text-based deep transformation functions that embed each data type separately and optimize unimodal embedding spaces such that they capture the topology of a disease KG or the text semantics of disease concepts. Finally, to achieve data fusion, REMAP aligns unimodal embedding spaces through a novel alignment penalty loss using shared disease concepts as anchors. This way, REMAP can effectively model data type-specific distribution and diverse representations while also aligning embeddings of distinct data types. Further, REMAP can be jointly trained on both graph and text data types but evaluated and implemented on either of the two modalities alone. In summary, the main contributions of this study are:

- We develop REMAP, a flexible multimodal approach for extracting and classifying disease-disease relations. REMAP fuses knowledge graph embeddings with deep language models and can flexibly accommodate missing data types, which is necessary to facilitate REMAP's validation and transition into biomedical implementation.

- We rigorously evaluate REMAP for extraction and classification of disease-disease relations. To this end, we create a training dataset using distant supervision and a high-quality test dataset of gold-standard annotations provided by three clinical domain experts. Evaluations show that REMAP achieves 88.6% micro-accuracy and 81.8% micro-F1 score on the human-annotated dataset, outperforming text-based methods by 10 and 17.2 percentage points, respectively. Further, REMAP achieves the best performance, 89.8% micro-accuracy, and 84.1% micro-F1 score, surpassing graph-based methods by 8.4 and 10.4 percentage points, respectively.

**JBI significance statement.** The significance of this study can be summarized as follows.

- **Problem or issue:** Enhance disease relation extraction through multimodal learning from language information and knowledge graphs.

- **What is already known:** Precise interpretation of disease-disease relationships is essential for building high-quality medical knowledge graphs. Although relation extraction based on either language information or knowledge graphs alone is a widely researched area, existing techniques are unable to realize the benefits arising from the confluence of language and graph information. Further, the applicability of multimodal methods is limited when data are incomplete or a subset of modalities is missing altogether.

- **What this paper adds:** We develop a flexible multi-modal approach for disease relation extraction that de-couples language and graphs but trains a joint model to address the challenge of

---

[2]Python implementation of REMAP is available on Github at `https://github.com/Lukeming-tsinghua/REMOD`. Our dataset of domain-expert annotations is at `https://doi.org/10.6084/m9.figshare.17776865`.

incomplete data modalities. We also built a human-annotated dataset and evaluate our method on it, demonstrating the effectiveness of the approach over both language- and graph-based methods.

## 2. Methods

We next detail the REMAP approach and illustrate it in detail for the task of disease relation extraction and classification (Figure 1). We first describe the notation, proceed with an overview of language and knowledge graph models, and outline the multimodal learning strategy to inject knowledge into extraction tasks.

### 2.1. Preliminaries

**Notation.** The input to REMAP is a combined dataset $D$ of language and graph information. This dataset consists of language information $D^T = \{B_i\}_{i=1}^{M_0} = D_L^T \cup D_U^T = \{B_i\}_{i=1}^{M} \cup \{B_i\}_{j=M+1}^{M_0}$ given as $M_0$ bags of sentences and graph information $D^G = \{(s_i, r_i, o_i)\}_{i=1}^{M} \bigcup \{(s_i, r_i, o_i)\}_{i=M+1}^{N}$ given as $N$ triplets $(s_i, r_i, o_i)$ encoding the relationship between $s_i$ and $o_i$ as $r_i$. For example, $s_i = $ "Hypobetalipoproteinemia", $o_i = $ "fatty liver" and $r_i = $ "May Cause" would indicate the fatty liver would be a possible symptom of hypobetalipoproteinemia. We assume that $M$ bags of sentences in $D^T$ overlap with the triplets from existing KG such that each sentence of the $i$th sentence bag contain $(s_i, r_i)$. The remaining $N - M$ triplets in $D^G$ cannot be mapped to sentences in $D^T$, and $M_0 - M$ sentences contain entity pairs that do not belong to existing KG. We represent the $i$-th sentence bag as $B_i = \{(t_{ij}, I_{ij}^s, I_{ij}^o)\}_{j=1}^{l_i}$, where $l_i$ is the number of sentences in bag $B_i$, $t_{ij}$ is the tokenized sequence of $j$-th sentence in $B_i$. Here, the tokenized sequence is a combination of the mentions of subject and object entities, entity markers, the document title, and the article structure. Marker tokens are added to each entity's head and tail position to denote entity type information. Last, $I_{ij}^s$ and $I_{ij}^o$ are start indices of entity markers for subject and object entities, respectively.

**Heterogeneous graph attention network.** Heterogeneous graph attention network (HAN) [30] is a graph neural network to embed a KG by leveraging meta paths. A meta path is a sequence of node types and relation types [58]. For example, in a disease KG, "Disease" → "May Cause" → "Disease" → "Differential Diagnosis" → "Disease" is a meta path. Node $u_i$ is connected to node $u_j$ via a meta path $\Phi$ if $u_i$ and $u_j$ are the head and tail nodes, respectively, of this meta path. Each node $u_i$ has an initial node embedding $\mathbf{h}_i^{\text{init}}$ and belongs to a type $\phi_i$, e.g., $\phi_i = $ "disease concept". Graph attention network specified a parameterized deep transformation function $f_{\text{HAN}}$ that maps nodes to condensed data summaries, i.e., embeddings, in a node-type specific manner as: $\mathbf{h}_i' = f_{\text{HAN}}^{\phi_i}(\mathbf{h}_i^{\text{init}})$.

We denote all nodes adjacent to $u_i$ via a meta-path $\Phi$ as $u_j \in N_i^{\Phi}$ and node-level attention mechanism provides information on how strongly $f_{\text{HAN}}$ attends to $u_i$'s each adjacent node $u_j$ when generating the embedding for $u_i$. In particular, the importance of $u_j$ for $u_i$ in meta path $\Phi$ is defined as:

$$a_{ij}^{\Phi} = \frac{\exp(\sigma(\mathbf{a}_{\Phi}^T \cdot [\mathbf{h}_i'||\mathbf{h}_j']))}{\sum_{k \in N_i^{\Phi}} \exp(\sigma(\mathbf{a}_{\Phi}^T \cdot [\mathbf{h}_i'||\mathbf{h}_k']))}, \tag{1}$$

where $\sigma$ is the sigmoid activation, $||$ indicates concatenation, and $\mathbf{a}_{\Phi}$ is a trainable vector. To promote stable attention, HAN uses multiple, i.e., $K$, heads and concatenates $K$ vectors after node level attention to produce the final node embedding for node $u_i$:

$$\mathbf{z}_i^{\Phi} = ||_{k=1}^{K} \sigma(\sum_{j \in N_i^{\Phi}} a_{ij}^{\Phi} \cdot \mathbf{h}_j'), \tag{2}$$

4

Given user-defined meta paths $\Phi_1, \ldots, \Phi_P$, HAN uses the above specified node-level attention to produce node embeddings $\mathbf{Z}_{\Phi_1}, \ldots, \mathbf{Z}_{\Phi_P}$. Finally, HAN uses semantic-level attention to combine meta path-specific node embeddings as:

$$\beta_{\Phi_p} = \frac{\exp(\frac{1}{|E|} \sum_{i \in E} \mathbf{q}^T \cdot \tanh(\mathbf{W} \cdot \mathbf{z}_i^{\Phi_p} + b))}{\sum_{p=1}^{P} \exp(\frac{1}{|E|} \sum_{i \in E} \mathbf{q}^T \cdot \tanh(\mathbf{W} \cdot \mathbf{z}_i^{\Phi_p} + b))}, \tag{3}$$

$$\mathbf{Z} = \sum_{p=1}^{P} \beta_{\Phi_p} \cdot \mathbf{Z}_{\Phi_p}, \tag{4}$$

where $\beta_{\Phi_p}$ represents the importance of meta path $\Phi_p$ towards final node embeddings $\mathbf{Z}$, and $\mathbf{q}^T$, $\mathbf{W}$, and $b$ are trainable parameters. The final outputs are node embeddings $\mathbf{Z}$, representing compact vector summaries of knowledge associated with each node in the KG.

**Translation- and tensor-based embeddings.** TransE [21] and TuckER [27] decode an optimized set of embeddings into the probability estimate of relationship $r_i$ existing between entities $s_i$ and $o_i$. This is achieved by a scoring function (SF) that either translates the embeddings in TransE as: $\mathrm{SF}_{\mathrm{Tr}}(\mathbf{h}_{s_i}, \mathbf{h}_{o_i}, \mathbf{h}_r) = \sigma(||\mathbf{h}_{s_i} + \mathbf{h}_r - \mathbf{h}_{o_i}||_2^2)$ or factorizes the embeddings in TuckER as: $\mathrm{SF}_{\mathrm{Tu}}(\mathbf{h}_{s_i}, \mathbf{h}_{o_i}, \mathbf{h}_r, \mathbf{W}) = \sigma(\mathbf{W} \times_1 \mathbf{h}_{s_i} \times_2 \mathbf{h}_r \times_3 \mathbf{h}_{o_i})$, where $\mathbf{h}_{s_i}$, $\mathbf{h}_{o_i}$, and $\mathbf{h}_r$ represent entity and relation embeddings, $\sigma$ denotes the sigmoid function, $\mathbf{W} \in \mathbb{R}^{d_{s_i} \times d_r \times d_{o_i}}$ is a trainable tensor, and $\times_d$, $d = 1, 2, 3$, indicates tensor multiplication along dimension $d$.

### 2.2. Text and knowledge graph encoders

**Embedding disease-associated sentences.** We start by tokenizing entities in sentences and proceed with an overview of the language encoder. Entity tokens identify the position and type of entities in a sentence [9, 10]. Specifically, tokens <S-type> and <S-type/>, and <O-type> and <O-type/> are used to denote the start and end of subject (S) and object (O) entities, respectively. The entity marker tokens are type-related, meaning that entities of different types (e.g., disease concepts, medications) get different tokens. This procedure produces bags of tokenized sentences $D^T$ that we encode into entity embeddings using a language encoder. We use SciBERT encoder [8] with the SciVocab vocabulary, which is a BERT language model optimized for scientific text with improved efficiency in biomedical domains than BioBERT or BERT model alone [10]. Tokenized sequences in a sentence bag $B_i$ are fed into the language model to produce a set of sequence outputs $\mathbb{H}_i = \{\mathbf{H}_i\}_{i=1}^{l_i}, i = 1, 2, \ldots, l_i$:

$$\mathbb{H}_i^{[m]} = \mathrm{SciBERT}(B_i; \mathbb{H}_i^{[m-1]}), \tag{5}$$

where $\mathbb{H}_i = [\mathbf{H}_{i1}, \ldots, \mathbf{H}_{il_i}]$, $l_i$ is the number of sentences in $B_i$, $\mathbf{H}_i \in \mathbb{R}^{d_l \times d_{hs}}$. We then aggregate representations of subject entities $s_i$ across all sentences in bag $B_i$ as: $\mathbf{H}_{s_i} = ||_{m=1}^{l_i} ||_{k \in I_{ij}^s} \mathbf{H}_{mk}$ and use self-attention to obtain the final language-based embedding $\mathbf{h}_{s_i}^T$ for subject entity $s_i$ as:

$$\mathbf{h}_{s_i}^T = \mathbf{H}_{s_i} \cdot \mathrm{softmax}(\boldsymbol{\omega} \cdot \tanh(\mathbf{H}_{s_i})), \tag{6}$$

where $\mathbf{h}_{s_i}^T \in \mathbb{R}^{d_{hs}}$ is the embedding of $s_i$ and $\boldsymbol{\omega}$ is a trainable vector. Embeddings of object entities (i.e., $\mathbf{h}_{o_i}^T$ for object entity $o_i$) are generated analogously by the language encoder. Self-attention is needed because specific sentences in a bag may not contain disease-disease relationships, and the attention mechanism allows the model to down-weight those uninformative sentences when generating embeddings.

**Embedding disease-disease knowledge relationships.** We use a heterogeneous graph attention encoder [30] to derive embeddings for nodes in the disease knowledge graph. The encoder produces embeddings for every subject entity $\mathbf{h}_{s_i}^G$ and every object entity $\mathbf{h}_{o_i}^G$ that have corresponding nodes in the KG as follows:

$$\mathbf{h}_{s_i}^G = f_{\text{HAN}}(D^G, \mathbf{H}^{\text{init}}, s_i), \quad \mathbf{h}_{o_i}^G = f_{\text{HAN}}(D^G, \mathbf{H}^{\text{init}}, o_i), \tag{7}$$

where $f_{\text{HAN}}$ transformation is given in Eqs. (1)-(4) and $\mathbf{H}^{\text{init}}$ denotes the matrix of initial embeddings.

**Scoring disease-disease relationships.** Taking language-based embeddings, $\mathbf{h}_{s_i}^T, \mathbf{h}_{o_i}^T$, and graph-based embeddings, $\mathbf{h}_{s_i}^G, \mathbf{h}_{o_i}^G$, for diseases that appear in either language or graph dataset, REMAP scores triplets $(s_i, r_k, o_i)$ as candidate disease-disease relationships. Specifically, to estimate the probability that diseases $s_i$ and $o_i$ are associated through a relation of type $r_k$ (e.g., $r_k =$ "May Cause"), REMAP calculates scores $p^T$ and $p^G$ representing the amount of evidence in the combined language-graph dataset that supports the disease-disease relationship:

$$p^T(r_{ik} = 1 | s_i, o_i) = \text{SF}(\mathbf{h}_{s_i}^T, \mathbf{h}_{o_i}^T, \mathbf{h}_r), \ k = 1, 2, \ldots, K, \tag{8}$$

$$p^G(r_{ik} = 1 | s_i, o_i) = \text{SF}(\mathbf{h}_{s_i}^G, \mathbf{h}_{o_i}^G, \mathbf{h}_r), \ k = 1, 2, \ldots, K, \tag{9}$$

where SF is the scoring function, and $K$ denotes the number of relation types. We consider three scoring functions, including the linear scoring function: $\text{SF}_{\text{Li}}(\mathbf{h}_{s_i}, \mathbf{h}_{o_i}, \mathbf{h}_{r_k}) = \sigma(\mathbf{W}_k(\mathbf{h}_{s_i} + \mathbf{h}_{o_i}) + b_k)$, the TransE scoring function: $\text{SF}_{\text{Tr}}(\mathbf{h}_{s_i}, \mathbf{h}_{o_i}, \mathbf{h}_{r_k}) = \sigma(||\mathbf{h}_{s_i} + \mathbf{h}_{r_k} - \mathbf{h}_{o_i}||_2^2)$, and the TuckER scoring function: $\text{SF}_{\text{Tu}}(\mathbf{h}_{s_i}, \mathbf{h}_{o_i}, \mathbf{h}_{r_k}, \mathbf{W}) = \sigma(\mathbf{W} \times_1 \mathbf{h}_{s_i} \times_2 \mathbf{h}_{r_k} \times_3 \mathbf{h}_{o_i})$, where separate kernels $\mathbf{W}^T$ and $\mathbf{W}^G$ are used for language and graph in the TuckER decomposition, and $\mathbf{h}_{r_k}$ denotes the encoding of relation type $r_k$ in TransE that is shared across both modalities (Section 2.1).

### 2.3. Co-training text and graph encoders

We proceed to describe the procedure for co-training text and graph encoders. From last section, we obtain relationship estimates based on evidence provided by text information, $p^T(r_{ik} = 1 | s_i, o_i)$, and graph information, $p^G(r_{ik} = 1 | s_i, o_i)$, for every triplet $(s_i, r_i, o_i)$. We use the binary cross entropy to optimize those estimates in each data type as:

$$L^T = \sum_{i=1}^{M} \sum_{k=1}^{K} r_{ik} \log(p^T(r_{ik} = 1 | s_i, o_i)) + (1 - r_{ik}) \log(1 - p^T(r_{ik} = 1 | s_i, o_i)), \tag{10}$$

$$L^G = \sum_{i=1}^{M} \sum_{k=1}^{K} r_{ik} \log(p^G(r_{ik} = 1 | s_i, o_i)) + (1 - r_{ik}) \log(1 - p^G(r_{ik} = 1 | s_i, o_i)), \tag{11}$$

and can combine language-based loss $L^T$ and graph-based loss $L^G$ as: $L_{\text{REMAP}} = L^T + L^G$. This loss function is motivated by the principle of knowledge distillation [59] to enhance multimodal interaction and improve classification performance. Using probabilities $p^T(r_{ik} = 1 | s_i, o_i)$ from the language encoder, we normalize them across $K$ relation types to obtain distribution $\mathbf{p}_t(s_i, o_i)$ as:

$$\mathbf{p}^T(s_i, o_i) = \{ \frac{e^{p^T(r_{ik} = 1 | s_i, o_i)}}{\sum_{m=1}^{K} e^{p^T(r_{im} = 1 | s_i, o_i)}} \}_{k=1}^{K}. \tag{12}$$

In the same manner, we calculate a graph-based distribution $\mathbf{p}^G(s_i, o_i)$ using softmax normalization.

Specifically, we develop two REMAP variants, REMAP-M and REMAP-B, based on how language-

based and graph-based losses are combined into a multimodal objective. In REMAP-M, both losses are aligned by shrinking the distance between distributions $\mathbf{p}^T$ and $\mathbf{p}^G$ using the Kullback-Leibler (KL) divergence:

$$D_{\text{KL}}(\mathbf{p}^T, \mathbf{p}^G) = \sum_{i=1}^{M} \sum_{k=1}^{K} p^G(s_i, o_i)_k \log(\frac{p^G(s_i, o_i)_k}{p^T(s_i, o_i)_k}), \tag{13}$$

where we measure the misalignment between language and graph models as follows:

$$L_{\text{REMAP-M}} = L_{\text{REMAP}} + \lambda_M(D_{\text{KL}}(\mathbf{p}^T, \mathbf{p}^G) + D_{\text{KL}}(\mathbf{p}^G, \mathbf{p}^T)). \tag{14}$$

Instead of measuring how distribution $\mathbf{p}^T$ is different from $\mathbf{p}^G$, REMAP-B selects the strongest logit across data types using an ensemble distillation strategy [60]. Specifically, REMAP-B uses the highest predicted score $p^B$ across both language and graph models to derive final predictions:

$$p^B(r_{ik} = 1|s_i, o_i) = \begin{cases} p^T(r_{ik} = 1|s_i, o_i), & p^T(r_{ik} = 1|s_i, o_i) \leq p^G(r_{ik} = 1|s_i, o_i) \text{ and } r_{ik} = 0 \\ p^T(r_{ik} = 1|s_i, o_i), & p^T(r_{ik} = 1|s_i, o_i) \geq p^G(r_{ik} = 1|s_i, o_i) \text{ and } r_{ik} = 1 \\ p^G(r_{ik} = 1|s_i, o_i), & p^T(r_{ik} = 1|s_i, o_i) \geq p^G(r_{ik} = 1|s_i, o_i) \text{ and } r_{ik} = 0 \\ p^G(r_{ik} = 1|s_i, o_i), & p^T(r_{ik} = 1|s_i, o_i) \leq p^G(r_{ik} = 1|s_i, o_i) \text{ and } r_{ik} = 1 \end{cases} \tag{15}$$

that are softmax-normalized across $K$ relation types:

$$\mathbf{p}^B(s_i, o_i) = \{\frac{e^{p^B(r_{ik}=1|s_i, o_i)}}{\sum_{m=1}^{K} e^{p^B(r_{im}=1|s_i, o_i)}}\}_{k=1}^{K}. \tag{16}$$

Finally, to minimize the discrepancy between predicted and known disease-disease relationships, REMAP-B uses minimizes the cross-entropy function:

$$L^B = \sum_{i=1}^{M} \sum_{k=1}^{K} r_{ik} \log(p^B(r_{ik} = 1|s_i, o_i)) + (1 - r_{ik}) \log(1 - p^B(r_{ik} = 1|s_i, o_i)), \tag{17}$$

and co-trains the multimodal encoder by promoting predictions that are aligned by both encoders:

$$L_{\text{REMAP-B}} = L_{\text{REMAP}} + \lambda_B[L^B + D_{\text{KL}}(\mathbf{p}^B, \mathbf{p}^T) + D_{\text{KL}}(\mathbf{p}^B, \mathbf{p}^G)]. \tag{18}$$

The outline of the complete REMAP-B algorithm is shown in Algorithm 1.

## 3. Experiments

We proceed with the description of datasets (Section 3.1), followed by implementation details of REMAP approach (Section 3.2) and the outline of experimental setup (Section 3.3).

### 3.1. Datasets

Datasets used in this study are multimodal and originate from diverse sources that we integrated and harmonized, as outlined below. In particular, we compiled a large disease-disease KG with text descriptions retrieved from medical data repositories using distant supervision following the data collection and preprocessing strategy described in Lin *et al.* [11]. Details on data preprocessing and feature engineering are described in Appendix A. Further, we utilize a large EHR dataset previously validated in Beam *et al.* [61] and a novel human-annotated dataset.

---

**Algorithm 1: REMAP, multimodal learning on graphs for disease relation extraction and classification.** Shown is the outline of REMAP-B (Section 2.3).

---

**Input:** Bag of sentences $\{B_i\}_{i=1}^n$, Knowledge graph $D^G = \{(s_i, r_i, o_i)\}$, Initial node embeddings $\mathbf{H}^{\text{init}}$, Scoring function SF, Regularization strength $\lambda_B$, Relation type vector pretrained on language dataset $\mathbf{r}^T$, Relation type vector pretrained on graph dataset $\mathbf{r}^G$

**Output:** Model parameters $\Theta$ of language-based and graph-based encoders

---

1   Initialize model parameters $\Theta$

2   Initialize relation representation as $\mathbf{r} = \frac{\mathbf{r}^T + \mathbf{r}^G}{2}$

3   **for** *epoch=1:n_epochs* **do**

4     **for** *i=1:n_bags* **do**

5       /* Encode language and the calculation logits, see Section 2.2 */

6       $\mathbf{H} = \text{SciBERT}(B_i)$

7       $\mathbf{H}_{s_i} = ||_{m=1}^{l_i} ||_{k \in I_{ij}^s} \mathbf{H}_{mk}$

8       $\mathbf{H}_{o_i} = ||_{m=1}^{l_i} ||_{k \in I_{ij}^o} \mathbf{H}_{mk}$

9       $\mathbf{h}_{s_i}^T = \mathbf{H}_{s_i} \cdot \text{softmax}(\boldsymbol{\omega} \cdot \tanh(\mathbf{H}_{s_i}))$

10      $\mathbf{h}_{o_i}^T = \mathbf{H}_{o_i} \cdot \text{softmax}(\boldsymbol{\omega} \cdot \tanh(\mathbf{H}_{o_i}))$

11      $\mathbf{p}^T(s_i, o_i) = \text{SF}(\mathbf{h}_{s_i}^T, \mathbf{h}_{o_i}^T, \mathbf{r})$

12      $L^T = \text{BinaryCrossEntropyLoss}(\mathbf{p}^T(s_i, o_i), \mathbf{r}(s_i, o_i))$

13      /* Encode graph and calculate graph logits, see Section 2.2 */

14      $\mathbf{h}_{s_i}^G = \text{HAN}(G, \mathbf{H}^{init}, s_i)$

15      $\mathbf{h}_{o_i}^G = \text{HAN}(G, \mathbf{H}^{init}, o_i)$

16      $p^G(s_i, o_i) = \text{SF}(\mathbf{h}_{s_i}^G, \mathbf{h}_{o_i}^G, \mathbf{r})$

17      $L^G = \text{BinaryCrossEntropyLoss}(\mathbf{p}^G(s_i, o_i), \mathbf{r}(s_i, o_i))$

18      /* Find the best logit and calculate alignment penalty loss, see Section 2.3 */

19      $\mathbf{p}^B(s_i, o_i) = \text{CalculateBestLogic}(\mathbf{p}^T(s_i, o_i), \mathbf{p}^G(s_i, o_i))$

20      $\mathbf{p}^T(s_i, o_i) = \text{softmax}(\mathbf{p}^T(s_i, o_i))$

21      $\mathbf{p}^G(s_i, o_i) = \text{softmax}(\mathbf{p}^G(s_i, o_i))$

22      $\mathbf{p}^B(s_i, o_i) = \text{softmax}(\mathbf{p}^B(s_i, o_i))$

23      $L^B = \sum_{i=1}^M \sum_{k=1}^K r_{ik} \log(p^B(r_{ik} = 1 | s_i, o_i)) + (1 - r_{ik}) \log(1 - p^B(r_{ik} = 1 | s_i, o_i))$

24      $L_{\text{REMAP-B}} = L^T + L^G + \lambda_B [L^B + D_{\text{KL}}(\mathbf{p}^B, \mathbf{p}^T) + D_{\text{KL}}(\mathbf{p}^B, \mathbf{p}^G)]$

25      $\Theta \leftarrow \text{Update}(\Theta, L_{\text{REMAP-B}})$

---

**Disease knowledge graph.** We construct a disease-disease KG from Diseases database [62] and MedScape [63] repository that unifies evidence on disease-disease associations based on text mining, manually curated literature, cancer mutation data, and genome-wide association studies. We construct a KG between diseases following the approach outlined in Lin *et al.* [11]. This KG contains 9,182 disease concepts, which are assigned concept unique identities (CUI) in the Unified Medical Language System (UMLS), and three types of relationships between disease concepts, which are 'may cause' (MC), 'may be caused by' (MBC), and 'differential diagnosis' (DDx). Other relations in the KG are denoted as 'not available' relations (NA). The MBC relation type is the reverse relation of the MC relation type, while DDx is a symmetric relation between disease concepts. Dataset statistics are in [11] and Table 1.

**Language dataset.** We use a text corpus taken from Lin *et al.* [11]. This corpus is built from medicine-related articles, including 42 million Pubmed abstracts[3], web pages collected via web

---

[3] Downloaded from The PubMed Baseline Repository provided by the National Library of Medicine, January 2021

crawler from 27 thousand pages on Uptodate and Medscape eMedicine, 10 thousand Wikipedia articles with medical titles, and the main text of four textbooks[4]. The corpus is segmented into sentences, which includes 237,119,572 sentences. We first employ forward maximum matching to identify all UMLS disease concepts in the corpus. Then, we link triplets in our disease knowledge graph to sentences if subject and object entities are both in the sentences. This method allows 1,466,065 sentences from these articles to be aligned to disease-disease edges in the disease knowledge graph. In the end, we group sentences by triplets. For example, the triplet *(Hypobetalipoproteinemia, may cause Fatty liver)* is aligned to a bag of sentences containing 4 sentences. More statistical details are shown in Table 1.

**Electronic health record dataset.** We use two types of information from electronic health records (EHRs), both taken from Beam *et al.* [61]. In particular, with a nationwide US health insurance plan with 60 million members over the period of 2008-2015, a dataset of concept co-occurrences from 20 million notes at Stanford, and an open access collection of 1.7 million full-text journal articles obtained from PubMed Central, Beam *et al.* [61] created a dataset of disease concepts identified with SNOMED-CT that appear together in the same note and used the SVD to decompose the resulting co-occurrence matrix and produce 500 dimension embedding vectors for disease concepts. We use the information on co-occurring disease concepts to guide the sampling of negative node pairs when training the disease knowledge graph. Further, we use the 500 dimension embedding vectors from [61] to initialize embeddings in REMAP's graph neural network. We use the average of all concept embedding as their initial embeddings for out-of-dictionary disease concepts.

**Human annotated dataset.** Relation triplets retrieved from databases may still have few incorrect data. To build a gold test set for a robust evaluation of our model, we randomly selected 500 disease triplets from our disease knowledge graph. We omit this subgraph from the knowledge graph used for model training and create an annotated dataset with it. For that, we recruited three senior physicians from Peking Union Medical College Hospital. We require them independently assign candidate relations *(May cause, May be caused, DDx and Not Available)* to these entity pairs without showing them relations collected from the database. We find annotation experts only disagreed on labels for 14 disease pairs (*i.e.*, 2.8% of the total number of disease pairs), and we resolve these disagreements through consensus voting. The human-annotated dataset is used for model comparison and performance evaluation.

*3.2. Training REMAP models*

Next, we outline the training details of REMAP models, including negative sampling, pre-training strategy for language and graph models, and the multimodal learning approach.

**Negative sampling.** We construct negative samples by sampling disease pairs whose co-occurrence in the EHR co-occurrence matrix [61] is lower than a pre-defined threshold. In particular, we expect that two unrelated diseases rarely appear together in EHRs, meaning that the corresponding values in the co-occurrence matrix are low. Thus, such disease pairs represent suitable negative samples to train models for classifying disease-disease relations.

**Pre-training a language model.** The text-based model comprises the text encoder and the TuckER module for disease-disease relation prediction. We denote the relation embeddings as $\mathbf{r}^T$, and the loss function as $L^T$ (Section 2.3). In particular, we use SciBERT tokenizer and SciBERT-SciVocab-uncased model [64]. The entity markers are added to the SciVocab vocabulary, and their

---

[4]*Harrison's Principles of Internal Medicine 20th Edition, Kelley's Textbook of Internal Medicine 4th Edition, Sabiston Textbook of Surgery: The Biological Basis of Modern Surgical Practice 20th Edition, and Kumar and Clark's Clinical Medicine 7th Edition*

**Table 1: Overview of the disease knowledge graph and the language dataset.** Shown are statistics for the following relation types: Not Available (NA), Differential Diagnosis (DDx), May Cause (MC), and May Be Caused by (MBC). Total denotes the total number of all triplets (see Figure 2).

| | Dataset | | Total | NA | DDx | MC | MBC | Entities |
|---|---|---|---|---|---|---|---|---|
| Unaligned | - | - | 96,913 | 30,546 | 20,657 | 23,411 | 22,298 | 9,182 |
| Aligned | Train | Triplet | 31,037 | 15,670 | 7,262 | 4,358 | 3,747 | 7,794 |
| | | Language | 1,244,874 | 799,194 | 208,921 | 123,735 | 113,024 | |
| | Validation | Triplet | 7,754 | 3,918 | 1,821 | 1,065 | 950 | 4,433 |
| | | Language | 206,179 | 68,934 | 60,165 | 43,706 | 33,474 | |
| | Annotated | Triplet | 499 | 8 | 210 | 159 | 122 | 733 |
| | | Language | 15,012 | 96 | 4,980 | 6,699 | 3,237 | |

embeddings are initialized with uniform distribution. We set the maximum number of sentences in a bag to $l_{m(max)}$. If the bag size is greater than $l_{m(max)}$, then $l_{m(max)}$ sentences are selected uniformly at random for training. Further details on hyper-parameters are in Table A3.

**Pre-training a graph neural network model.** The graph-based model comprises the heterogeneous attention network encoder and the TuckER module for disease-disease relation prediction. The relation embeddings produced by the TuckER module are denoted as $\mathbf{r}^G$. In the pre-training phase, the model is optimized for the loss function is $L^G$ (Section 2.3). The initial embeddings for nodes in the disease knowledge graph are concept unique identifier (CUI) representations derived from the SVD decomposition of the EHR co-occurrence matrix [61]. Further details on hyper-parameters are in Table A3.

**Cross-modal learning.** After data type-specific pre-training is completed, the text and graph models are fused in cross-modal learning. To this end, the shared relation vector $\mathbf{r}$ is initialized as: $\mathbf{r} = (\mathbf{r}^T + \mathbf{r}^G)/2$. We consider two REMAP variants, namely REMAP-M and REMAP-B, optimized for different loss functions (Section 2.3). Details on hyper-parameter selection are in Table A3.

*3.3. Experimental setup*

Next we overview baseline methods and performance metrics.

**Baseline methods.** We consider 9 baseline methods divided into two groups: 5 methods for link prediction on KGs and 5 methods for relation extraction in text.

For graph-based baselines, we take both disease and relation embeddings as parameters. We randomly initialize the embeddings and train the knowledge graph models or graph neural networks. Graph-based baselines are trained on the disease KG and include the following:

- **TransE** [21] embeds diseases and relations by translating embedding vectors in the learned embedding space. Given embeddings of a triplet $(\mathbf{h}, \mathbf{r}, \mathbf{t})$, the score is calculated as $s = ||\mathbf{h} + \mathbf{r} - \mathbf{t}||_2^2$. We train our TransE model using negative sampling and $L2$ penalty as recommended by the authors. The negative samples are constructed by randomly replacing objects $\mathbf{t}$ given subject $\mathbf{h}$ and relation $\mathbf{r}$. We begin with random embeddings for diseases and relations and then optimize them with the same margin-based ranking criterion and stochastic gradient descent as in the original TransE paper.

- **DistMult** [65] predicts disease-disease relationships using a bilinear decoder for edge in the KG. Given embeddings of a triplet $(\mathbf{h}, \mathbf{r}, \mathbf{t})$, the score is calculated as $s = \mathbf{h}\mathbf{M}_r\mathbf{t}$, where $\mathbf{M}_r$ is a

trainable diagonal matrix for relation $r$. The disease and relation embeddings are parameters in this baseline. We initialize it with random embeddings and optimize them with margin-based ranking loss and batch stochastic gradient descent. The negative samples are constructed by corrupting the subject or object in a relation triplet.

- **ComplEx** [66] predicts disease-disease relationships by carrying out a matrix factorization using complex-valued embeddings. Denoting the complex-valued embeddings of a relation triplet as $\mathbf{h}_s$, $\mathbf{w}_r$, and $\mathbf{h}_o$. ComplEx assumes $P(Y = 1) = \sigma(\phi(\mathbf{h}_s, \mathbf{w}_r, \mathbf{h}_o))$, where $Y = 1$ represents the triplet $(h, r, t)$ holds and $\sigma$ denotes the sigmoid function. The score function $\phi(\cdot)$ is calculated as $s = Re(< \mathbf{h}, \mathbf{r}, \mathbf{t} >)$. The product $< \cdot >$ is a Hermitian product, and relation embedding $\mathbf{r}$ is a complex-valued vector. We take random embeddings as initial embeddings for both real and imaginary parts. The training object is minimizing the negative log-likelihood of this logistic model with $L_2$ penalty on the parameters of disease and relation embeddings.

- **RGCN** [67] predicts disease-disease relationships using relational graph convolutional network (RGCN). Following the authors' recommendations, we use a 2-layer RGCN to embed the KG and DistMult decoder for link prediction. We train the RGCN model using the cross-entropy loss on four relations and the Adam optimizer. The RGCN uses the sigmoid activation function and has a 0.4 dropout rate for self-loops and 0.2 for other edges.

- **TuckER** [33] is a KG embedding method that uses Tucker tensor decomposition. Given the embeddings of a triplet $(\mathbf{h}, \mathbf{r}, \mathbf{t})$, the score is calculated as $s = \mathbf{W} \times_1 \mathbf{h} \times_2 \mathbf{r} \times_3 \mathbf{t}$, where $\mathbf{W}$ is a trainable matrix and $\times_i$ denotes tensor multiplication on dimension $i$.

Text-based methods are trained on the language dataset following Lin *et al.* [11] strategy. We consider models with bag-of-words (BoW) engineered features, convolutional neural networks, and pre-trained language models. Sentence-level baselines, including RF, TextCNN, BiGRU, and PubmedBERT, use majority voting on sentence-level predictions to produce final predictions for input examples that are longer than one sentence. We consider the following text-based baselines:

- **Random Forest (RF)** uses scaled BoW features with 100 decision trees and the Gini criterion as the classifier. Punctuations and English stop words are removed from sentences [68]. We extract 40,000 most frequently occurring N-grams and calculate the TF-IDF to produce sentence-level features to RF. The N-gram TF-IDF transformation is a widely-used and effective feature construction procedure.

- **TextCNN** [69] is a convolutional neural network for text, a useful deep learning algorithm for sentence classification tasks, such as sentiment analysis and question classification. The TextCNN we take as a baseline has 64 feature maps for each size and 8 different sizes of feature maps whose lengths of filter windows range from 2 to 10. We use ReLU as the activation function, 1 dimension max pooling after convolutional layers, and a dropout layer with a 0.25 dropout rate. We use skip-gram embeddings [70] to initialize the TextCNN model and train it using cross-entropy loss and Adam optimizer.

- **BiGRU** [71] is a recurrent neural network initialized in the same way as TextCNN. The encoder is a 1-layer BiGRU whose output hidden states are aggregated into sentence embedding using a one-headed self-attention. These sentence embeddings are fed into a linear layer and provide logits for relation classification. The dimension of input word embeddings is 200, and the dimension of BiGRU hidden states is 100. We also add dropout layers with a 0.5 dropout

rate after BiGRU and the linear layer. The sentence-level predictions given by this model will generate an instance-level classification result by majority voting.

- **BiGRU+Attention** [71] is the same as the BiGRU baseline with an added instance-level attention. Instead of voting, this approach uses one-head instance-level attention to aggregate latent vectors across all sentences in an instance. This provides an instance embedding for each sentence bag. Then we use a linear layer to provide logits of relation classification from instance embedding. This baseline is trained with cross-entropy loss on four relations, including Not Available (NA). We use the same hyper-parameters for this baseline as for BiGRU.

- **PubmedBERT** [72] is a domain-specific pre-trained language model trained on the corpus of abstracts from Pubmed with mask language modeling and next sentence prediction as pretext tasks. We take the checkpoint of PubmedBERT and fine tune the model for disease relation extraction. We use the output embedding of the [CLS] token as sentence embedding. And then sentence embedding in a bag is aggregated with instance-level attention. The fine tuning employs the Adam optimizer and a linear scheduler to adjust the learning rate. The loss function is also cross-entropy loss on four different relations, including Not Available (NA).

We use grid search on the validation set to select hyper-parameters for all methods. Table A3 outlines the hyper-parameter selection in REMAP.

**Performance metrics.** We evaluate predicted disease-disease relations by calculating the accuracy of predicted relations between disease concepts, which is an established approach for benchmarking relation extraction methods [73]. Specifically, given a triplet $(s_i, r_i, o_i)$ and a predicted score $p_i \in [0, 1]$, relation $r_i$ is predicted to exist between $s_i$ and $o_i$ if the predicted score $p_i \geq \text{threshold}_i$, which corresponds to a binary classification task for each relation type. The $\text{threshold}_i$ is a relation type-specific value determined such that binary classification performance achieves maximal F1-score. We report classification accuracy, precision, recall, and F1-score for all experiments in this study.

## 3.4. Variants of REMAP approach

We carry out an ablation study to examine the utility of key REMAP components. We consider the following three components and examine REMAP's performance with and without each. The results of all variants are reported in the ablation study 3.

- **REMAP-B without joint learning** In text-only ablations, we use SciBERT to obtain concept embeddings $\mathbf{h}_s^T$ and $\mathbf{h}_o^T$, and combine them to produce a relation embedding $\mathbf{r}_k$, where $\mathbf{r}_k$ is related to the concept embeddings based on text information. Finally, we consider three scoring functions to classify disease-disease relations, and we denote the models as SciBERT (linear), SciBERT (TransE), and SciBERT (TuckER). Similarly, in graph-only ablations, we first use a heterogeneous attention network to obtain graph embeddings $\mathbf{h}_s^G$ and $\mathbf{h}_o^G$ that are combined into prediction by different scoring functions, including HAN (linear), HAN (TransE), and HAN (TuckER).

- **REMAP-B without EHR embeddings** REMAP-B uses EHR embeddings [61] as initial node embedding $\mathbf{H}^{init}$. To examine the utility of EHR embeddings, we design an ablation study that initializes node embeddings using the popular Xavier initialization [74] instead of EHR embeddings. Other parts of the model are the same as in REMAP-B.

- **REMAP-B without unaligned triplets** Unaligned triplets denote triplets in the disease knowledge graph that do not have the corresponding sentences in the language dataset. To

demonstrate how these unaligned triplets influence model performance, we design an ablation study in which we train a REMAP-B model on the reduced disease knowledge graph with unaligned triplets excluded.

## 4. Results

REMAP is a multimodal language-graph learning approach. We evaluate REMAP's prowess for disease relation extraction when the REMAP model is tasked to identify candidate disease relations in either text (Section 4.1) or graph-structured (Section 4.2) data. We present ablation study and case studies in the discussion (Section 5).

### 4.1. Extracting disease relations from text

We start with results on the human-annotated dataset where each method, while it can be trained on a multimodal text-graph dataset, is tasked to extract disease relations from text alone, meaning that the test set consists of only annotated sentences. Table 2 shows performance on the annotated set for text-based methods. We witness the pre-trained language model, such as PubmedBERT[72], consistently outperforms other neural network baselines and random forest in F1 score, which achieves 78.5 in micro average. Further, BiGRU+Attention is the best performing baseline method in accuracy, achieving an accuracy of 78.6. We also find that REMAP-B and REMAP-M achieve the best performance across all settings, outperforming baselines significantly. In particular, REMAP models surpass the strongest baseline by 10.0 absolute percentage points (accuracy) and by 7.2 absolute percentage points (F1-score). These results show that multimodal learning can considerably advance disease relation extraction when only one data type is available at test time.

### 4.2. Completing disease KG with novel disease-disease relationships

We proceed with the results of disease relation prediction in a setting where each method is tasked to classify disease relations based on the disease knowledge graph alone. This setting evaluates the flexibility of REMAP as REMAP can answer either graph-based or text-based disease queries. In particular, Table 2 shows performance results attained on the human-annotated dataset with query disease pairs given only as disease concept nodes in the knowledge graph. We found TuckER significantly outperforms other knowledge graph embedding baselines in both accuracy and F1 score. Last, we find that REMAP-B is a top performer among REMAP variants, achieving an accuracy of 89.8 and an F1-score of 84.1.

## 5. Discussion

Next, we analyze results focusing on the selection of negative disease-disease samples and the impact negative sampling has on model performance. We also examine trade-offs between translation and bilinear methods (Section 5.1). Finally, we give a case study illustrating REMAP predictions (Section 5.2) and provide an ablation study into key REMAP components (Section 5.3).

### 5.1. Further analysis of results

Our dataset's negative samples are more representative than the previous approaches [11] of generating negative samples. These negative samples are selected from the EHR co-occurrence matrix, which has the co-occurrence number of disease-disease pairs below the threshold we set. In this paper, the negative samples constructed in this way do not require conceptual replacement (that is, they are not "generated" negative samples), so the negative samples are more in line with

**Table 2: Results of disease relation extraction on the human-annotated set.** DDx: differential diagnosis, MC: may cause, MBC: may be caused by. The "micro" columns denote micro average accuracy or F1-score for DDx, MC, and MBC relation types. Further results are in .

| Modality | | Model | Accuracy | | | | F1-score | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | micro | DDx | MC | MBC | micro | DDx | MC | MBC |
| Text | Baselines | RF | 72.0 | 68.6 | 70.8 | 76.4 | 37.8 | 53.1 | 25.5 | 19.2 |
| | | TextCNN | 76.7 | 75.4 | 73.0 | 81.8 | 60.9 | 67.5 | 59.9 | 48.6 |
| | | BiGRU | 77.4 | 73.0 | 77.2 | 82.0 | 62.0 | 67.9 | 54.0 | 59.8 |
| | | BiGRU+attention | 78.6 | 75.0 | 78.2 | 82.6 | 64.6 | 67.7 | 63.5 | 60.6 |
| | | PubmedBERT | 77.7 | 82.4 | 75.5 | 82.0 | 78.5 | 75.2 | 74.6 | 81.7 |
| | Ours | REMAP | 88.2 | 83.6 | 89.0 | 92.0 | 80.9 | 80.7 | 80.0 | 82.6 |
| | | REMAP-M | 88.6 | 84.2 | 89.0 | **92.8** | 81.5 | 81.6 | 79.6 | **83.8** |
| | | REMAP-B | **88.6** | **84.4** | **89.2** | 92.4 | **81.8** | **81.9** | **80.3** | 83.3 |
| Graph | Baselines | TransE_l2 | 75.1 | 70.7 | 72.7 | 81.8 | 63.2 | 68.0 | 57.0 | 62.2 |
| | | DistMult | 69.8 | 77.5 | 61.3 | 70.5 | 56.1 | 71.0 | 43.4 | 51.5 |
| | | ComplEx | 79.0 | 75.2 | 77.8 | 84.2 | 65.0 | 69.3 | 56.5 | 66.9 |
| | | RGCN | 71.8 | 78.6 | 62.5 | 74.3 | 62.2 | 75.1 | 50.8 | 58.6 |
| | | TuckER | 81.5 | 77.6 | 82.3 | 84.7 | 73.7 | 76.2 | 71.7 | 71.9 |
| | Ours | REMAP | 89.6 | 86.4 | 89.6 | **92.8** | 83.5 | 84.3 | 81.6 | **84.2** |
| | | REMAP-M | 89.3 | 87.0 | 88.4 | 92.6 | 83.3 | 85.7 | 78.8 | 83.8 |
| | | REMAP-B | **89.8** | **87.3** | **89.9** | 92.2 | **84.1** | **85.8** | **82.4** | 82.7 |

the actual situation. At the same time, threshold control is also used to reduce the false negative rate. However, the negative samples constructed in this way make it more difficult to distinguish between traditional machine learning and deep learning models, so a more fusion-capable model is needed for mining.

Further, we find that TuckER obtained better performance than TransE in the experiments. We hypothesize that this finding is due to the inherent limitation of TransE in the presence of 1-to-N relationships. Specifically, suppose multiple target diseases exist for a particular source disease and a relation. In that case, the representation returned by TransE fails to capture the local graph neighborhoods of all target diseases simultaneously. That is because $\mathbf{t} \approx \mathbf{h} + \mathbf{r}$, meaning that the model does not have sufficient capacity to differentiate between embeddings $\mathbf{t}$ for different target diseases. In contrast, bilinear representation of TuckER can address this limitation of TransE, which explains its better performance. For this reason, our REMAP approach uses the TuckER component to facilitate joint learning.

### 5.2. Case study

We proceed with a case study examining the May Cause (MC) relationship between hypobetalipoproteinemia and fatty liver disease to illustrate how REMAP uses both graph and text modalities to classify the MC type disease-disease relationship. Several studies [75, 76] found the hypobetalipoproteinemia can cause fatty liver. Figure 3 illustrates the prediction of the MC relationship between hypobetalipoproteinemia and fatty liver made by the REMAP-B joint learning model. The graph can provide graph structure information to REMAP-B. For example, hypobetalipoproteinemia has 8 outgoing edges of MC type, which are the most among all edge types, and fatty liver disease

**Table 3: Results of the ablation study.** DDx: differential diagnosis, MC: may cause, MBC: may be caused by. The "micro" columns denote micro average accuracy or F1-score for DDx, MC, and MBC relation types.

| Modality | Model | Accuracy | | | | F1-score | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | micro | DDx | MC | MBC | micro | DDx | MC | MBC |
| Text | REMAP-B | **88.6** | **84.4** | **89.2** | 92.4 | **81.8** | **81.9** | 80.3 | 83.3 |
| | w/o joint learning (linear) | 87.3 | 84.1 | 86.3 | 91.3 | 78.9 | 81.2 | 73.4 | 80.9 |
| | w/o joint learning (TransE) | 86.1 | 79.1 | 89.1 | 90.2 | 80.5 | 79.3 | **82.2** | 80.8 |
| | w/o joint learning (TuckER) | 87.9 | 83.6 | 87.0 | **93.2** | 80.0 | 80.1 | 75.7 | 85.1 |
| | w/o EHR embedding | 87.6 | 83.2 | 88.2 | 91.4 | 79.6 | 79.8 | 78.2 | 81.1 |
| | w/o unaligned triplets | 88.2 | 83.0 | 88.5 | **93.2** | 81.0 | 80.0 | 78.8 | **85.2** |
| Graph | REMAP-B | **89.8** | **87.3** | **89.9** | **92.2** | **84.1** | **85.8** | **82.4** | **82.7** |
| | w/o joint learning (linear) | 87.4 | 82.4 | 89.3 | 90.2 | 80.5 | 80.5 | 82.1 | 78.4 |
| | w/o joint learning (TransE) | 85.8 | 81.6 | 85.3 | 90.4 | 77.2 | 78.9 | 73.1 | 78.9 |
| | w/o joint learning (TuckER) | 88.9 | 85.8 | 89.4 | 91.6 | 82.5 | 83.9 | 81.0 | 81.4 |
| | w/o EHR embedding | 87.6 | 84.4 | 87.2 | 91.4 | 80.3 | 82.0 | 76.8 | 81.4 |
| | w/o unaligned triplets | 87.3 | 84.8 | 87.8 | 89.4 | 79.5 | 82.4 | 77.7 | 76.2 |

also has 4 outgoing edges of the May Be Caused (MBC) type. And the graph encoder can also capture information from any possible meta-paths linking these two nodes in a shape of Hypobetalipoproteinemia $\underrightarrow{MC}$ X $\underrightarrow{MC}$ Fatty liver, where $X$ is also a node in the knowledge graph. The language encoder can capture language information from the free text, semantic types, and spans of disease from special tokens we added to the vocabulary of the pretrained language model. In the case of joint learning, the text-based model can extract part of the disease representation from the knowledge graph to update its internal representations and thus improve text-based classification of relations.

### 5.3. Ablation study

We conduct an ablation study to provide evidence for the effectiveness of four key components in REMAP-B: joint learning loss, scoring functions, EHR embedding, and unaligned triplets in the knowledge graph. Table 3 shows results that compare performance REMAP-B and performance after excluding each component. We observe performance drops when excluding joint learning techniques in REMAP-B. This conclusion is consistent across almost all the metrics and edge types in both text and graph modality, except the edge type MBC (May Be Caused) in text modality, on which the accuracy increases 0.8 from 92.4 to 93.2 percent when excluding joint learning with TuckER. We also find that EHR embedding plays a significant role in REMAP-B since the performance consistently decreases in all metrics, edge types, and both modalities. The largest drop in performance is observed on the F1-score in graph modality, which decreases by 3.8 percent from 84.1 to 80.3. This analysis demonstrates that REMAP-B can effectively employ the EHR data for disease relation extraction.

Further, accuracy and F1-score also decrease when we remove unaligned triplets in the knowledge graph in all settings except the edge types MBC (May Be Caused), on which the accuracy increases 0.8 percent and F1-score increases 1.9 percent. This result indicates the importance of leveraging unsupervised information from unaligned triplets. In summary, all components we introduce in REMAP-B are vital for achieving outstanding performance on disease relation extraction. Among all these components, joint learning techniques are most important since performance drops considerably

when we exclude them, and this result is consistent across all settings.

## 6. Conclusion

We develop a multi-modal learning approach REMAP for disease relation extraction that fuses language modeling with knowledge graphs. Results on a dataset of clinical expert annotations show that REMAP considerably outperforms methods for learning on text or knowledge graphs alone. Further, REMAP can extract and classify disease relationships in the most challenging settings where text or graph information is absent. Finally, we provide a new data resource of extracted relationships between diseases that can serve as a benchmarking dataset for systematic evaluation and comparison of disease relation extraction algorithms.

## Funding and acknowledgments

## Author contributions

Y.L and K.L. are co-first authors and have developed the method and carried out all analyses in this study. S.Y. contributed experimental data. T.C. and M.Z. conceived and designed the study. The project was guided by M.Z., including designing methodology and outlining experiments. Y.L., K.L., S.Y., T.C., and M.Z. wrote the manuscript. All authors discussed the results and reviewed the paper.

## Data and code availability

Python implementation of REMAP is available on Github at https://github.com/Lukeming-tsinghua/REMOD. The human annotated dataset used for evaluation of REMAP is available at https://doi.org/10.6084/m9.figshare.17776865.

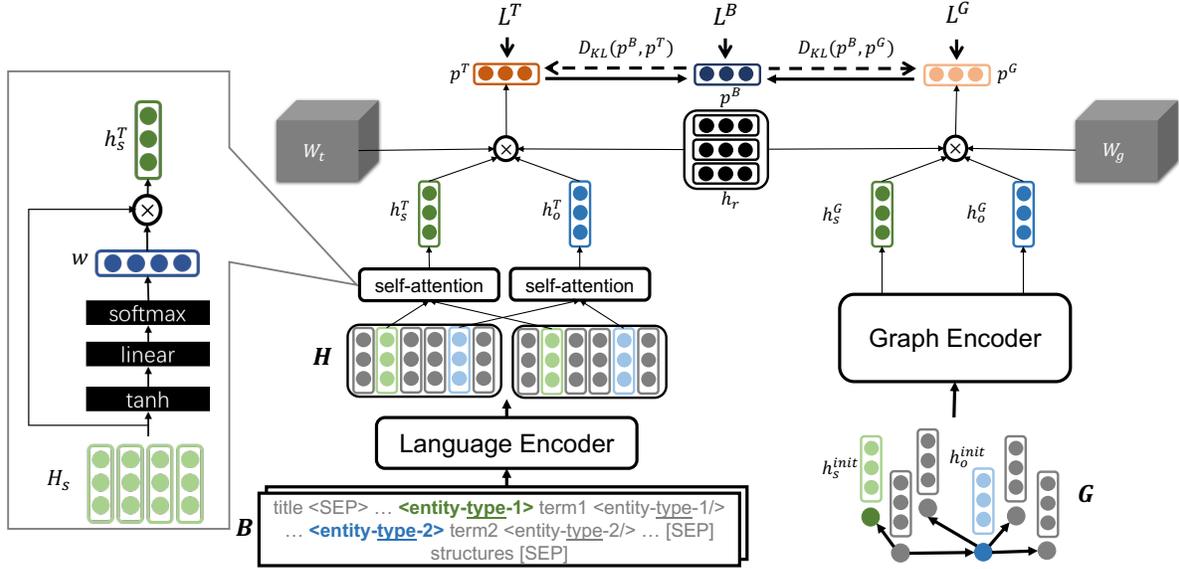## Conflicts of interest

None declared.

**Figure 1: Overview of REMAP architecture.** REMAP introduces a novel co-training learning strategy that continually updates a multimodal language-graph model for disease relation extraction and classification. Language and graph encoders specify deep transformation functions that embed disease concepts (i.e., subject entities $s$ and object entities $o$) from the language data $D^T$ and disease knowledge graph $D^G$ into compact embeddings, producing condensed summaries of language semantics and biomedical knowledge for every disease. Embeddings output by the encoders (i.e., $\mathbf{h}_s^T$, $\mathbf{h}_o^T$, $\mathbf{h}_o^G$, $\mathbf{h}_s^G$) are then combined in a disease relation type-specific manner (e.g., "differential diagnosis" and "may cause" relation types) and passed to a scoring function that calculates the probability representing how likely two diseases are related to each other and what kind of relationship exists between them.
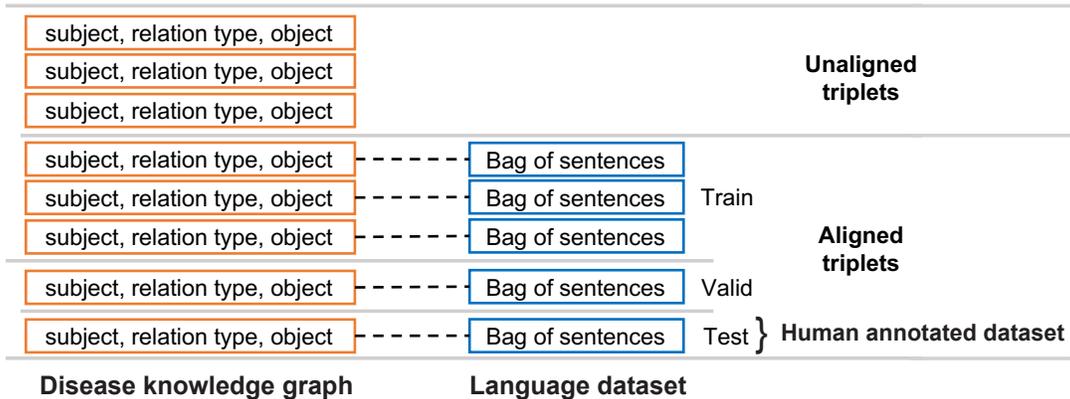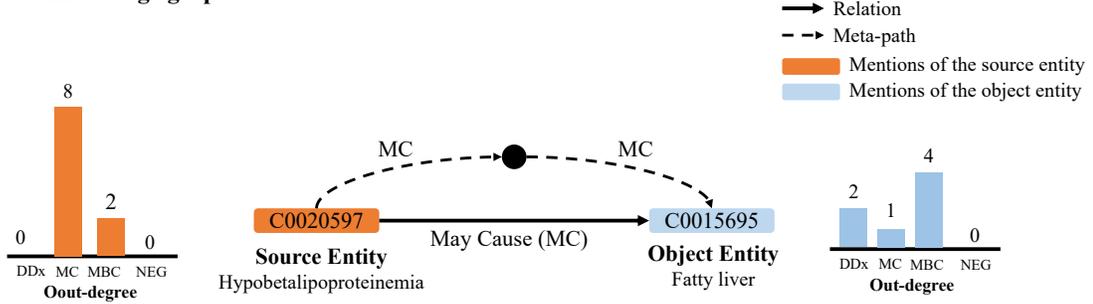


**Figure 2: Creating a dataset split for performance evaluation and benchmarking.** Triplets in the disease knowledge graph are divided into aligned triplets and unaligned triplets. *Aligned triplets* are triplets that have corresponding sentences in the language dataset. *Unaligned triplets* have no corresponding sentences in the language dataset.

**a) Disease knowledge graph**

8

2

0    0
DDx MC MBC NEG
**Oout-degree**

→ Relation
--→ Meta-path
■ Mentions of the source entity
■ Mentions of the object entity

MC          MC

C0020597 ——— May Cause (MC) ———→ C0015695
**Source Entity**                **Object Entity**
Hypobetalipoproteinemia          Fatty liver

4

2
1
0
DDx MC MBC NEG
**Out-degree**

**b) Medical text**

1. **\<empty_title\> \<sep\>** humans and genetically engineered mice with **\<entity-t047-1\>** hypobetalipoproteinemia **\</entity-t047-1\>** due to truncation-producing mutations of the apolipoprotein b ( apob ) gene frequently have \<entity-t047-2\> fatty liver **\</entity-t047-2\>** , because the apob defect impairs the capacity of livers to export triglycerides ( tgs ) .

2. **\<empty_title\> \<sep\>** based on these rare findings , heterozygous **\<entity-t047-1\>** hypobetalipoproteinemia **\</entity-t047-1\>** should thus be considered as a possible cause in patients presenting with an unexplained **\<entity-t047-2\>** fatty liver **\</entity-t047-2\>** .

3. **\<empty_title\> \<sep\>** heterozygous **\<entity-t047-1\>** hypobetalipoproteinemia **\</entity-t047-1\>** should be considered in a hypolipidemic subject with an otherwise unexplained **\<entity-t047-2\>** fatty liver **\</entity-t047-2\>** .

**Figure 3: Illustration of (Hypobetalipoproteinemia, May Cause, Fatty liver) triplet in REMAP.** This triplet represents the Hypobetalipoproteinemia has the symptom of fatty liver. The subject entity (Hypobetalipoproteinemia, C0020597) is shown in orange and the object entity (Fatty liver, C0015695) is shown in blue. **(a)** The subject and object entities are identified with unique UMLS concept identifiers. The relation between them is the may cause (MC). There is a MC-MC meta-path between these entities and that information is leveraged by our graph encoder to predict the relationship between Hypobetalipoproteinemia and Fatty liver. The bar plots indicate distribution of relation types going out of the subject or object entities in the knowledge graph. **(b)** We show three sentences representing the triplet. We use $<sep>$ token to separate the sentences and the title of article from which we mine the sentence. If the article is from PubMed, we use special token $<empty\_title>$ as a placeholder. We add two special tokens before and after terms in each sentence to identify the positions of entities. For example, we add $<entity\text{-}t047\text{-}1>$ before *Hypobetalipoproteinemia* to mark the beginning of subject entity; *t047* serves as a identifier of the semantic type - *Disease or Syndrome* - for Hypobetalipoproteinemia.

## Appendix  A. Data preprocessing

Data preprocessing and feature engineering follow the strategy utlined in Lin *et al.* [11]. We outline the data preparation process in this section. This process transforms the raw datasets, including relation triplets, text corpus, and electronic health records embeddings into the AI-ready data for AI analyses. The code is publicly available on GitHub[5].

### *Appendix  A.1. Relation triplet collection*

Acquiring relation triplets is the first step in the preparation of training data. The majority of disease relation triplets are directly collected from the Diseases Database [62]. We also collect relation triplets by resolving semi-structured content on the MedScape. Pages on it usually follow a very standard template, which allows one to easily locate sections about the target relations, where the entities are usually presented in lists and tables. Most commonly, the page title provides the head entity, the section title specifies the relation, lists, and tables in the section provide the tail entities. Therefore, we write simple web scraping scripts and apply maximum mapping to identify mentions of entities with UMLS CUIs, and assemble them into relation triplets.

### *Appendix  A.2. Document preparation*

The input to this step is a set of documents collected from online sources. These documents are used to pretrain a model *en_core_web_sm* in *SpaCy* to split the document into sentences with additional structure information, including the title and subject headings.

### *Appendix  A.3. Entity linking*

The input to entity linking are sentences extracted from corpus. We use the nested forward maximum matching algorithm to annotate the mentions with medical terms collected from UMLS.

### *Appendix  A.4. Data cleaning*

We discard sentences that are shorter than 5 words. And we generate the static position embedding for baseline models, such as TextCNN and RNNs. We also make sure the head entity appears before the tail entity in sentences, and those sentences with entities in reverse order are moved to the opposite relation (i.e., from may_cause to may_be_caused_by).

## Appendix  B. Further information on REMAP's performance

Table A1 extends Table 2 from the main text with additional information about REMAP's performance measured using precision, recall, and F1-score metrics. Shown is performance that REMAP achieves when the REMAP model, although trained on the multimodal graph-text dataset, is asked to make predictions at test time using only text information.

Table A2 extends Table 2 from the main text with additional information about REMAP's performance measured using precision, recall, and F1-score metrics. Shown is performance that REMAP achieves when the REMAP model, although trained on the multimodal graph-text dataset, is asked to make predictions at test time using only information from the disease knowledge graph.

## Appendix  C. Hyper parameters

---

[5]https://github.com/lychyzclc/High-throughput-relation-extraction-algorithm

**Table A1: Performance of multimodal methods for identifying candidate disease-disease relations <u>from text data.</u>** Shown is average performance across multiple independent runs calculated on the human annotated set. Higher values indicate better performance.

| Model | Precision | | | | Recall | | | | F1-score | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | minor | DDx | MC | MBC | minor | DDx | MC | MBC | minor | DDx | MC | MBC |
| Random Forest | 69.2 | 71.8 | 67.6 | 58.3 | 26.0 | 42.2 | 15.7 | 58.3 | 37.8 | 53.1 | 25.5 | 19.2 |
| TextCNN | 67.8 | 76.2 | 56.7 | 78.2 | 55.3 | 60.7 | 63.5 | 35.2 | 60.9 | 67.5 | 59.9 | 48.6 |
| BiGRU | 69.1 | 68.1 | 75.3 | 65.7 | 56.3 | 67.8 | 42.1 | 54.9 | 62.0 | 67.9 | 54.0 | 59.8 |
| BiGRU+attention | 70.6 | 74.4 | 67.9 | 67.7 | 59.6 | 62.1 | 59.7 | 54.9 | 64.6 | 67.7 | 63.5 | 60.6 |
| PubmedBERT | 71.2 | 64.0 | 71.2 | 78.4 | 87.4 | 91.2 | 85.5 | 85.3 | 78.5 | 75.2 | 74.6 | 81.7 |
| SciBERT (linear) | 86.6 | 81.0 | 96.9 | 88.3 | 72.5 | 81.4 | 59.1 | 74.6 | 78.9 | 81.2 | 73.4 | 80.9 |
| SciBERT (TransE) | 74.9 | 68.2 | 86.2 | 77.4 | 87.0 | 94.8 | 78.6 | 84.4 | 80.5 | 79.3 | 82.2 | 80.8 |
| SciBERT (TuckER) | 87.3 | **81.7** | 93.5 | 91.5 | 73.9 | 78.6 | 63.5 | **79.5** | 80.0 | 80.1 | 75.7 | **85.1** |
| REMAP | 85.8 | 79.9 | 94.8 | 88.0 | 76.6 | 81.4 | **69.2** | 77.9 | 80.9 | 80.7 | 80.0 | 82.6 |
| REMAP-M | **87.4** | 79.9 | **97.3** | **93.0** | 76.4 | 83.3 | 67.3 | 76.2 | 81.5 | 81.6 | 79.6 | 83.8 |
| REMAP-B | 86.2 | 79.7 | 95.7 | 89.6 | **77.8** | **84.3** | **69.2** | 77.9 | **81.8** | **81.9** | **80.3** | 83.3 |

**Table A2: Performance of multimodal methods for identifying candidate disease-disease relations <u>from graph-structured data.</u>** Shown is average performance across multiple independent runs calculated on the human annotated set. Higher values indicate better performance.

| Model | Precision | | | | Recall | | | | F1-score | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | minor | DDx | MC | MBC | minor | DDx | MC | MBC | minor | DDx | MC | MBC |
| TransE_l2 | 61.3 | 63.0 | 57.3 | 63.0 | 65.2 | 73.8 | 56.6 | 61.5 | 63.2 | 68.0 | 57.0 | 62.2 |
| DistMult | 53.6 | 77.8 | 40.7 | 43.1 | 58.9 | 65.2 | 46.5 | 63.9 | 56.1 | 71.0 | 43.4 | 51.5 |
| ComplEx | 71.7 | 72.2 | 75.0 | 68.4 | 59.5 | 66.7 | 45.3 | 65.6 | 65.0 | 69.3 | 56.5 | 66.9 |
| RGCN | 56.2 | 74.9 | 44.2 | 48.9 | 69.7 | 75.2 | 59.7 | 73.0 | 62.2 | 75.1 | 50.8 | 58.6 |
| TuckER | 70.1 | 69.8 | 74.3 | 66.2 | 77.8 | 83.8 | 69.2 | 78.7 | 73.7 | 76.2 | 71.7 | 71.9 |
| HAN (linear) | 82.0 | 75.2 | 92.9 | 84.8 | 79.0 | 86.7 | **73.6** | 73.0 | 80.5 | 80.5 | **82.1** | 78.4 |
| HAN (TransE) | 81.3 | 76.1 | 88.4 | 84.9 | 73.5 | 81.9 | 62.3 | 73.8 | 77.2 | 78.9 | 73.1 | 78.9 |
| HAN (TuckER) | 85.7 | 80.1 | 94.2 | 88.5 | 79.4 | 88.1 | 71.1 | 75.4 | 82.5 | 83.9 | 81.0 | 81.4 |
| REMAP | **87.0** | **81.7** | 93.5 | **90.6** | 80.2 | 87.1 | 72.3 | **78.7** | 83.5 | 84.3 | 81.6 | **84.2** |
| REMAP-M | 85.6 | 79.8 | **93.9** | 89.7 | 81.1 | **92.4** | 67.9 | **78.7** | 83.3 | 85.7 | 78.8 | 83.8 |
| REMAP-B | 86.6 | 81.3 | 93.6 | 90.3 | **81.7** | 91.0 | **73.6** | 76.2 | **84.1** | **85.8** | **82.4** | 82.7 |

**Table A3: Selection of hyper-parameters in REMAP**. We use grid search to select hyper-parameters on the validation set.

| REMAP's component | Model parameter | Notation | Value |
|---|---|---|---|
| Neural architecture | Padding length of sentences | $d_l$ | 256 |
| | Hidden size of SciBERT output | $d_{hs}$ | 768 |
| | Hidden size of HAN output | $d_{ha}$ | 100 |
| | Hidden size of initial node embedding | $d_{hi}$ | 1,000 |
| | Hidden size of node embedding | $d_h$ | 100 |
| | Hidden size of relation embedding | $d_r$ | 100 |
| Multimodal training | Max sentence sample number | $l_{m(\max)}$ | 12 |
| | Training batch size | $b_{\text{train}}$ | 4 |
| | Test batch size | $b_{\text{test}}$ | 16 |
| | Weight decay | $wd$ | $5 \times 10^{-5}$ |
| | Learning rate | $lr$ | $1 \times 10^{-5}$ |
| | Gradient accumulate step | $\text{step}_g$ | 4 |
| | Optimizer | | Adam |
| | Scheduler | | Linear |
| | Warmup rate | $r_{\text{warmup}}$ | 0.1 |
| Language model | Max sentence sample number | $l_{m(\max)}$ | 12 |
| | Training batch size | $b_{\text{train}}$ | 4 |
| | Test batch size | $b_{\text{test}}$ | 16 |
| | Weight decay | $wd$ | $5 \times 10^{-5}$ |
| | Learning rate | $lr$ | $1 \times 10^{-5}$ |
| | Gradient accumulate step | $\text{step}_g$ | 4 |
| | Optimizer | | Adam |
| | Scheduler | | Linear |
| | Warmup rate | $r_{\text{warmup}}$ | 0.1 |
| Graph model | Training batch size | $b_{\text{train}}$ | 512 |
| | Test batch size | $b_{\text{test}}$ | 512 |
| | Weight decay | $wd$ | $1 \times 10^{-8}$ |
| | Learning rate | $lr$ | $1 \times 10^{-3}$ |
| | Optimizer | | Adam |
| | Scheduler | | StepLR |
| | StepLR scheduler | $\gamma$ | 0.9 |

# References

[1] C. Ruiz, M. Zitnik, J. Leskovec, Identification of disease treatment mechanisms through the multiscale interactome, Nature Communications 12 (1) (2021) 1–15.

[2] C. Hong, E. Rush, M. Liu, D. Zhou, J. Sun, A. Sonabend, V. M. Castro, P. Schubert, V. A. Panickan, T. Cai, et al., Clinical knowledge extraction via sparse embedding regression (KESER) with multi-center large scale electronic health record data, NPJ Digital Medicine 4 (1) (2021) 1–11.

[3] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, J. Taylor, Freebase: a collaboratively created graph database for structuring human knowledge, in: ACM SIGMOD, 2008, pp. 1247–1250.

[4] T. C. Rindflesch, H. Kilicoglu, M. Fiszman, G. Rosemblat, D. Shin, Semantic medline: An advanced information management application for biomedicine, Information Services & Use (2011) 15–21.

[5] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka, T. M. Mitchell, Toward an architecture for never-ending language learning, in: AAAI, 2010.

[6] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmann, S. Sun, W. Zhang, Knowledge vault: A web-scale approach to probabilistic knowledge fusion, in: ACM SIGKDD, 2014, pp. 601–610.

[7] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: NAACL-HLT, 2019, pp. 4171 – 4186.

[8] I. Beltagy, K. Lo, A. Cohan, SciBERT: A pretrained language model for scientific text, in: EMNLP-IJCNLP, 2019, pp. 3615–3620.

[9] L. B. Soares, N. FitzGerald, J. Ling, T. Kwiatkowski, Matching the blanks: Distributional similarity for relation learning, in: ACL, 2019, pp. 2895–2905.

[10] Z. Zhong, D. Chen, A frustratingly easy approach for joint entity and relation extraction, in: NAACL-HT, 2021, pp. 50–61.

[11] Y. Lin, K. Lu, Y. Chen, C. Hong, S. Yu, High-throughput relation extraction algorithm development associating knowledge articles and electronic health records, arXiv:2009.03506 (2020).

[12] J. Chen, B. Hu, W. Peng, Q. Chen, B. Tang, Biomedical relation extraction via knowledge-enhanced reading comprehension, BMC Bioinformatics (2022) 1–19.

[13] G. Li, C. Wu, K. Vijay-Shanker, Noise reduction methods for distantly supervised biomedical relation extraction, in: BioNLP, 2017, pp. 184–193.

[14] Q. Wang, Z. Mao, B. Wang, L. Guo, Knowledge graph embedding: A survey of approaches and applications, IEEE Transactions on Knowledge and Data Engineering 29 (12) (2017) 2724–2743.

[15] M. M. Li, K. Huang, M. Zitnik, Representation learning for networks in biology and medicine: advancements, challenges, and opportunities, arXiv:2104.04883 (2021).

[16] M. Zitnik, B. Zupan, Collective pairwise classification for multi-way analysis of disease and drug data, in: the Pacific Symposium on Biocomputing, 2016, pp. 81–92.

[17] B. Shi, T. Weninger, ProjE: Embedding projection for knowledge graph completion, in: AAAI, Vol. 31, 2017.

[18] Y. Lin, Z. Liu, M. Sun, Y. Liu, X. Zhu, Learning entity and relation embeddings for knowledge graph completion, in: AAAI, 2015.

[19] X. Wang, X. He, Y. Cao, M. Liu, T.-S. Chua, KGAT: Knowledge graph attention network for recommendation, in: KDD, 2019, pp. 950–958.

[20] Z. Sun, J. Yang, J. Zhang, A. Bozzon, L.-K. Huang, C. Xu, Recurrent knowledge graph embedding for effective recommendation, in: ACM Conference on Recommender Systems, 2018, pp. 297–305.

[21] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, O. Yakhnenko, Translating embeddings for modeling multi-relational data, Advances in Neural Information Processing Systems 26 (2013) 2787–2795.

[22] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph and text jointly embedding, in: EMNLP, 2014, pp. 1591–1601.

[23] G. Ji, S. He, L. Xu, K. Liu, J. Zhao, Knowledge graph embedding via dynamic mapping matrix, in: ACL, 2015, pp. 687–696.

[24] M. Nickel, V. Tresp, H.-P. Kriegel, A three-way model for collective learning on multi-relational data, in: ICML, 2011.

[25] B. Yang, W. tau Yih, X. He, J. Gao, L. Deng, Embedding entities and relations for learning and inference in knowledge bases, Computing Research Repository abs/1412.6575 (2015).

[26] T. Trouillon, C. R. Dance, J. Welbl, S. Riedel, É. Gaussier, G. Bouchard, Knowledge graph completion via complex tensor factorization, JMLR 18 (2017).

[27] I. Balazevic, C. Allen, T. Hospedales, TuckER: Tensor factorization for knowledge graph completion, in: EMNLP-IJCNLP, 2019, pp. 5185–5194.

[28] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. Van Den Berg, I. Titov, M. Welling, Modeling relational data with graph convolutional networks, in: European Semantic Web Conference, Springer, 2018, pp. 593–607.

[29] D. Busbridge, D. Sherburn, P. Cavallo, N. Y. Hammerla, Relational graph attention networks, arXiv:1904.05811 (2019).

[30] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, P. S. Yu, Heterogeneous graph attention network, in: WWW, 2019, pp. 2022–2032.

[31] M. Zitnik, M. Agrawal, J. Leskovec, Modeling polypharmacy side effects with graph convolutional networks, Bioinformatics 34 (13) (2018) i457–i466.

[32] M. Mintz, S. Bills, R. Snow, D. Jurafsky, Distant supervision for relation extraction without labeled data, in: ACL-IJCNLP, 2009, pp. 1003–1011.

[33] I. Balažević, C. Allen, T. M. Hospedales, Tucker: Tensor factorization for knowledge graph completion, in: EMNLP, 2019.

[34] W. Liu, P. Zhou, Z. Zhao, Z. Wang, Q. Ju, H. Deng, P. Wang, K-BERT: Enabling language representation with knowledge graph, in: AAAI, Vol. 34, 2020, pp. 2901–2908.

[35] Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, H. Wu, ERNIE: Enhanced representation through knowledge integration, arXiv:1904.09223 (2019).

[36] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, Q. Liu, ERNIE: Enhanced language representation with informative entities, in: ACL, 2019, pp. 1441–1451.

[37] B. He, D. Zhou, J. Xiao, X. Jiang, Q. Liu, N. J. Yuan, T. Xu, BERT-MK: Integrating graph contextualized knowledge into pre-trained language models, in: Findings of EMNLP, 2020, pp. 2281–2290.

[38] R. Koncel-Kedziorski, D. Bekal, Y. Luan, M. Lapata, H. Hajishirzi, Text generation from knowledge graphs with graph transformers, arXiv:1904.02342 (2019).

[39] T. Sun, Y. Shao, X. Qiu, Q. Guo, Y. Hu, X. Huang, Z. Zhang, Colake: Contextualized language and knowledge embedding, in: International Conference on Computational Linguistics, 2020, pp. 3660–3670.

[40] L. Hu, L. Zhang, C. Shi, L. Nie, W. Guan, C. Yang, Improving distantly-supervised relation extraction with joint label embedding, in: EMNLP, 2019, pp. 3812–3820.

[41] P. Xu, D. Barbosa, Connecting language and knowledge with heterogeneous representations for neural relation extraction, NAACL-HLT (2019) 3201–3206.

[42] N. Zhang, S. Deng, Z. Sun, G. Wang, X. Chen, W. Zhang, H. Chen, Long-tail relation extraction via knowledge graph embeddings and graph convolution networks, in: NAACL-HLT, 2019, pp. 3016–3025.

[43] Y. Wang, H. Zhang, G. Shi, Z. Liu, Q. Zhou, A model of text-enhanced knowledge graph representation learning with mutual attention, IEEE Access 8 (2020) 52895–52905.

[44] X. Han, Z. Liu, M. Sun, Joint representation learning of text and knowledge for knowledge graph completion, arXiv:1611.04125 (2016).

[45] Z. Ji, Z. Lei, T. Shen, J. Zhang, Joint representations of knowledge graphs and textual information via reference sentences, IEICE Transactions on Information and Systems (2020) 1362–1370.

[46] Q. Dai, N. Inoue, P. Reisert, R. Takahashi, K. Inui, Distantly supervised biomedical knowledge acquisition via knowledge graph based attention, in: the Workshop on Extracting Structured Knowledge from Scientific Publications, 2019, pp. 1–10.

[47] G. Stoica, E. A. Platanios, B. Póczos, Improving relation extraction by leveraging knowledge graph link prediction, arXiv:2012.04812 (2020).

[48] M. Zitnik, F. Nguyen, B. Wang, J. Leskovec, A. Goldenberg, M. M. Hoffman, Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities, Information Fusion 50 (2019) 71–91.

[49] Z. He, W. Chen, Y. Wang, W. Zhang, G. Wang, M. Zhang, Improving neural relation extraction with positive and unlabeled learning, in: AAAI, Vol. 34, 2020, pp. 7927–7934.

[50] P. H. Le-Khac, G. Healy, A. F. Smeaton, Contrastive representation learning: A framework and review, IEEE Access (2020).

[51] P. Su, Y. Peng, K. Vijay-Shanker, Improving bert model using contrastive learning for biomedical relation extraction, in: BIONLP, 2021.

[52] Q. Wang, L. Zhan, P. Thompson, J. Zhou, Multimodal learning with incomplete modalities by knowledge distillation, in: SIGKDD, 2020, p. 1828–1838.

[53] Q. Suo, W. Zhong, F. Ma, Y. Yuan, J. Gao, A. Zhang, Metric learning on healthcare data with incomplete modalities., in: IJCAI, 2019, pp. 3534–3540.

[54] T. Zhou, M. Liu, K.-H. Thung, D. Shen, Latent representation learning for alzheimer's disease diagnosis with incomplete multi-modality neuroimaging and genetic data, IEEE Transactions on Medical Imaging (2019) 2411–2422.

[55] Y. Yang, D.-C. Zhan, X.-R. Sheng, Y. Jiang, Semi-supervised multi-modal learning with incomplete modalities., in: IJCAI, 2018, pp. 2998–3004.

[56] L. Cai, Z. Wang, H. Gao, D. Shen, S. Ji, Deep adversarial learning for multi-modality missing data completion, in: ACM SIGKDD, 2018, pp. 1158–1166.

[57] N. Jaques, S. Taylor, A. Sano, R. Picard, Multimodal autoencoder: A deep learning approach to filling in missing sensor data and enabling better mood prediction, in: ACII, 2017, pp. 202–208.

[58] Y. Sun, J. Han, X. Yan, P. S. Yu, T. Wu, PathSim: meta path-based top-k similarity search in heterogeneous information networks, in: VLDB, 2011, pp. 992–1003.

[59] X. Lan, X. Zhu, S. Gong, Knowledge distillation by on-the-fly native ensemble, in: NeurIPS, 2018.

[60] Q. Guo, X. Wang, Y. Wu, Z. Yu, D. Liang, X. Hu, P. Luo, Online knowledge distillation via collaborative learning, in: CVPR, 2020, pp. 11020–11029.

[61] A. L. Beam, B. Kompa, A. Schmaltz, I. Fried, G. Weber, N. Palmer, X. Shi, T. Cai, I. S. Kohane, Clinical concept embeddings learned from massive sources of multimodal medical data, in: Pacific Symposium on Biocomputing 2020, 2019, pp. 295–306.

[62] Laptop diseases database ver 2.0; medical lists and links diseases database [internet]. [cited 2020 aug 16]. available from: http://www.diseasesdatabase.com/.

[63] P. Frishauf, Medscape–the first 5 years, Medscape General Medicine (2005) 5.

[64] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., Huggingface's transformers: State-of-the-art natural language processing, arXiv:1910.03771 (2019).

[65] B. Yang, W. tau Yih, X. He, J. Gao, L. Deng, Embedding entities and relations for learning and inference in knowledge bases, ICLR (2015).

[66] T. Trouillon, J. Welbl, S. Riedel, E. Gaussier, G. Bouchard, Complex embeddings for simple link prediction, ICML (2016).

[67] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. v. d. Berg, I. Titov, M. Welling, Modeling relational data with graph convolutional networks, arXiv:1703.06103 (2017).

[68] C. D. Manning, M. Surdeanu, J. Bauer, J. R. Finkel, S. Bethard, D. McClosky, The stanford corenlp natural language processing toolkit, in: Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations, 2014, pp. 55–60.

[69] Y. Kim, Convolutional neural networks for sentence classification, in: EMNLP, 2014, pp. 1746–1751.

[70] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, Advances in neural information processing systems 26 (2013).

[71] K. Cho, B. Van Merriënboer, D. Bahdanau, Y. Bengio, On the properties of neural machine translation: Encoder-decoder approaches, arXiv preprint arXiv:1409.1259 (2014).

[72] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, H. Poon, Domain-specific language model pretraining for biomedical natural language processing, ACM Transactions on Computing for Healthcare (HEALTH) (2021) 1–23.

[73] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, J. Leskovec, Open Graph Benchmark: Datasets for machine learning on graphs, NeurIPS (2020).

[74] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in: AISTATS, 2010, pp. 249–256.

[75] G. Schonfeld, B. W. Patterson, D. A. Yablonskiy, T. S. Tanoli, M. Averna, N. Elias, P. Yue, J. Ackerman, Fatty liver in familial hypobetalipoproteinemia: triglyceride assembly into vldl particles is affected by the extent of hepatic steatosis, Journal of lipid research 44 (3) (2003) 470–478.

[76] J. Rodrigues, A. Azevedo, S. Tavares, C. Rocha, E. S. Silva, Non-alcoholic fatty liver disease associated with hypobetalipoproteinemia: report of three cases and a novel mutation in apob gene, NASCER E CRESCER-BIRTH AND GROWTH MEDICAL JOURNAL 25 (2) (2016) 104–107.