

# DYNAMICALLY PRUNED MESSAGE PASSING NETWORKS FOR LARGE-SCALE KNOWLEDGE GRAPH REASONING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We propose *Dynamically Pruned Message Passing Networks* (DPMPN) for large-scale knowledge graph reasoning. In contrast to existing models, embedding-based or path-based, we learn an input-dependent subgraph to explicitly model a sequential reasoning process. Each subgraph is dynamically constructed, expanding itself selectively under a flow-style attention mechanism. In this way, we can not only construct graphical explanations to interpret prediction, but also prune message passing in Graph Neural Networks (GNNs) to scale with the size of graphs. We take the inspiration from the consciousness prior proposed by Bengio (2017) to design a two-GNN framework to encode global input-invariant graph-structured representation and learn local input-dependent one coordinated by an attention module. Experiments show the reasoning capability in our model that is providing a clear graphical explanation as well as predicting results accurately, outperforming most state-of-the-art methods in knowledge base completion tasks.

## 1 INTRODUCTION

Modern deep learning systems need to acquire the reasoning capability beyond their black-box nature to produce interpretable predictions (Pearl & Mackenzie, 2018; Bengio, 2018). In what form we model a reasoning process should be given more thought than just obtaining a final prediction. Intuitively, a reasoning process can be regarded as a sequence of using existing facts to establish new knowledge, step by step, and finally drawing conclusions in the form of constructing explanations as well as making predictions. Therefore, it needs an explicit modeling to identify and organize reasoning steps to form a clear interpretable representation during predicting. A natural idea is to use graph-structured representation where a semantic unit or pairwise relation can be explicitly represented by a node or an edge as building blocks to support graph-based reasoning, a more flexible form in contrast to rigid deductive logical reasoning (Battaglia et al., 2018; Xu et al., 2019).

Graph-based reasoning can be applied to a wide variety of real-world scenarios. Here, we choose knowledge graph-related tasks to explore due to its representativeness. In knowledge base completion (KBC) tasks, embedding-based models (Bordes et al., 2013; Yang et al., 2015; Dettmers et al., 2018; Trouillon et al., 2016; Sun et al., 2018; Lacroix et al., 2018) can easily obtain a very competitive score by fitting data using various neural network techniques, but lacking an explicit modeling to construct explanations by directly exploiting graph structure prevents it from being interpretable, a critical property of reasoning, since Euclidean embedding space will not produce a clearly stated and human-readable representation.

Recent work for knowledge graph (KG) reasoning focuses on path-based (Wang, 2018; Xiong et al., 2017; Das et al., 2018; Shen et al., 2018; Chen et al., 2018; Lin et al., 2018) or logic-like models (Cohen, 2016; Yang et al., 2017). Most of them construct an explicit path to model an iterative decision-making process using reinforcement learning and recurrent networks. However, a question is: do we have a better form, more flexible and interpretable, to express reasoning in the graph context rather than one or several paths. To this end, we propose to learn a subgraph starting from a head node and expanding itself conditionally and selectively according to a query relation, where a tail node is predicted after the last expansion. To better explain how the tail is determined by the expansion, we weigh, prune and save intermediate nodes selected at each step to capture long-range

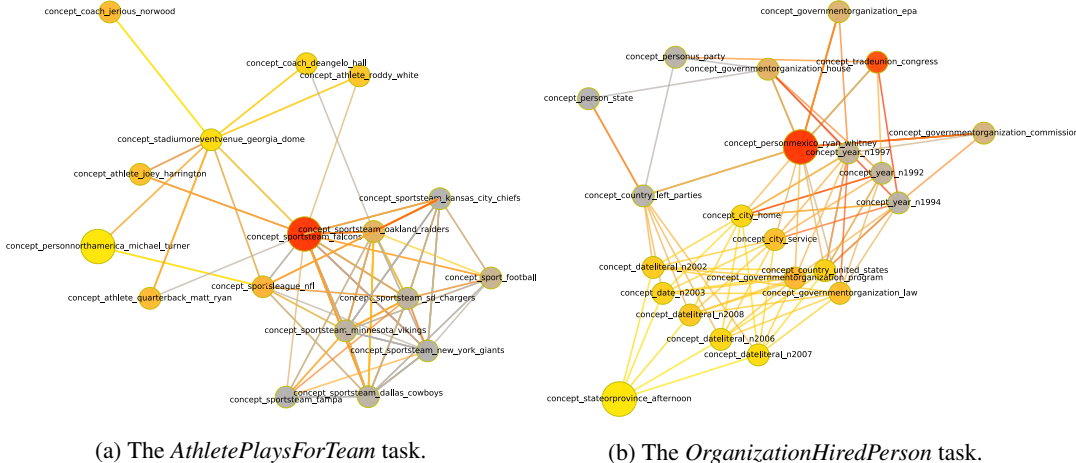


Figure 1: Subgraph visualization of reasoning results on two examples from NELL995’s test data. One is for the *AthletePlaysForTeam* task and the other for the *OrganizationHiredPerson* task. Each task has a graph with ten thousands of nodes and edges. The big yellow in each part represents a given head and the big red represents a predicted tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

dependence and yield a concise and compact subgraph explanation for the tail prediction as shown in Figure 1.

**Graph reasoning can be powered by Graph Neural Networks.** Graph reasoning demands a way to effectively learn about entities, relations, and rules for composing them, that is, an ability for combinatorial generalization by manipulating structured knowledge and producing structured explanations. Graph Neural Networks (GNNs) provide such structured representation and computation and also inherit powerful data-fitting capacity from deep neural networks (Scarselli et al., 2009; Battaglia et al., 2018). Specifically, GNNs follow a neighborhood aggregation scheme, recursively aggregating and transforming neighboring nodes’ representations to update representations for each node. Therefore, after  $T$  iterations of aggregation, each node can carry the structured information within the node’s  $T$ -hop neighborhood (Gilmer et al., 2017; Xu et al., 2018a).

**GNNs need graphical attention expression to interpret.** Neighborhood attention operation is a popular way to implement attention mechanism on graphs (Velickovic et al., 2018; Hoshen, 2017) by using multi-head self-attention to focus on specific interactions with neighbors when aggregating messages. However, we argue that graphical attention expression should be designed instead not only to facilitate structured computation but also to construct dynamically pruned structured explanations. We present three considerations: (1) selecting nodes based on currently operated subgraphs, that is, first attending over nodes within subgraphs to pick a smaller set and then attending over the picked nodes’ neighbors to expand subgraphs, (2) breaking isolation of attention operations used for each step and propagating attention across steps like a flow to produce long-range influence, and (3) that such flow-style attention mechanism should model a changing node probability distribution, that is, a Markov process driven by step-varying transition matrices. Besides, we should use an attention module disentangled from representation aggregating and transforming in GNNs to explicitly model a reasoning process on graphs out of low-level representation computing.

**GNNs need input-dependent pruning to scale.** GNNs are notorious for its poor scalability due to its heavy computation complexity. Consider, for example, one message passing iteration performed over a graph with  $|V|$  nodes and  $|E|$  edges. It has quadratic complexity in the number of nodes,  $O(|V|^2)$ , if the graph is fully connected. Even if the graph is sparse so that the complexity can be reduced to  $O(|E|)$  by exploiting structural sparsity, it is still problematic when meeting large graphs with millions of nodes and edges. Besides, mini-batch based training with batch size  $B$  and high dimensions  $D$  would make things worse, leading to the complexity of  $O(BD|E|)$ . We argue that this situation can be avoided by learning input-dependent pruning, as in most cases an input example uses a small fraction of the entire graph, and it is wasteful to perform homogeneous structured computation over the full graph for each input. Therefore, we propose to prune message passing depending on inputs and run on dynamical computation graphs instead of a static computation graph.

**Cognitive intuition of the consciousness prior.** The notion of attentive awareness has been shared by cognitive science communities in several theories (Dehaene et al., 1998; Tononi et al., 2016). Bengio (2017) brought this notion into deep learning models in his *consciousness prior* proposal. He pointed out a process of disentangling high-level abstract factors from full underlying representation to form a low-dimensional combination of a few selected factors to constitute a conscious thought, and emphasized the role of attention in expressing awareness during this process. Bengio proposed to use two recurrent neural networks (RNNs) to encode two types of state: the unconscious state represented by a full high-dimensional vector before applying attention, and the conscious state by a derived low-dimensional vector after applying attention.

In our work, we use two GNNs to encode such states into node representation vectors. However, standard message passing runs globally so that messages gathered by a node can come from everywhere and get further entangled by aggregation operations. Therefore, we draw an input-dependent or context-aware local subgraph to constrain message passing. We also want to access global information about the graph structure to get a boarder view before focusing on a local subgraph. Inspired by the consciousness prior, we apply attention mechanism to the two GNNs, where the bottom one performs input-invariant standard message passing globally, called **Inattentive GNN (IGNN)**, and the above one performs input-dependent pruned message passing locally, called **Attentive GNN (AGNN)**. The intuition is that the Inattentive GNN can support the Attentive GNN by providing raw representation, entangled but rich, while the Attentive GNN captures various input-dependent subgraphs consisting of a few selected nodes and their edges, cohesive with sharp semantics, disentangled from the full graph. Nodes within such a subgraph are more densely connected to form a small community to further exchange information and make decisions collectively on how to grow the subgraph next. In experiments, we find our model can run on very large graphs with millions of edges, such as the YAGO3-10 dataset, even using a laptop without causing out-of-memory errors. Our prediction results of KBC tasks attain very competitive scores on HITS@1,3 and the mean reciprocal rank (MRR) compared to the best embedding-based method so far, and we provide interpretations that they do not have.

## 2 PROBLEM FORMULATION

**Notation.** We use a supervised setting with training data  $f(x_i, y_i)g_{i=1}^N$  where  $x_i$  is an input and  $y_i$  is a target. We denote a full graph by  $G$  with node set  $V$  and edge set  $E$ , and denote an input-dependent subgraph by  $G(x)$  with node set  $V_{G(x)}$  and edge set  $E_{G(x)}$ . We also denote the set of edge types (or relation types) by  $R$ . We require each subgraph  $G$  to hold  $E_G = f(v, u) \geq E : v, u \geq V_G g$ , so that we can define  $G = G_1 [ G_2$  if  $V_G = V_{G_1} [ V_{G_2}$  and define  $G_1 \sim G_2$  if  $V_{G_1} \sim V_{G_2}$ . We define the boundary of a subgraph as  $\partial G$  if  $V_{\partial G} = N(V_G) \sim V_G$  where  $N(V_G)$  means the union of neighbors of all the nodes in  $V_G$ . We also define high-order boundaries such as  $\partial^2 G$  if  $V_{\partial^2 G} = N(N(V_G)) [ N(V_G) \sim V_G$ . Trainable parameters include node embeddings  $f_{e_v}g_{v \geq 2}V$ , relation type embeddings  $f_{e_r}g_{r \geq 2}R$ , and neural network weights used in two GNNs and an attention module. When performing standard or pruned message passing, node embeddings and relation type embeddings will be indexed according to the operated graph, and thus we denote them by  $\theta_G$  or  $\theta_{G(x)}$ . We denote batch size by  $B$  and dimensions by  $D$ . For IGNN, we use  $\mathcal{H}^t$  of size  $jVj \cdot D$  to denote node hidden states at step  $t$ ; for AGNN, we use  $\mathbf{H}^t(x)$  of size  $jV_{G(x)}j \cdot D$  to denote.

We define the objective based on our two GNNs as  $\sum_{i=1}^N l(x_i, y_i; \theta_{G(x_i)}, \theta_G)$ , where  $G(x_i)$  is dynamically constructed. First, we write the standard message passing in IGNN as

$$\mathcal{H}^t = f_{\text{IGNN}}(\mathcal{H}^{t-1}; \theta_G), \quad (1)$$

where  $f_{\text{IGNN}}$  represents all involved operations in one message passing iteration over  $G$ , including: (1) computing messages along each edge with the complexity<sup>1</sup> of  $O(BDjEj)$ , (2) aggregating messages received at each node with the complexity of  $O(BDjEj)$ , and (3) updating node hidden states with the complexity of  $O(BDjVj)$ . For a  $T$ -step propagation, we get the per-batch complexity of  $O(BDT(jEj + jVj))$ . Considering that backpropagation requires intermediate computation results to be saved during one pass, this complexity counts for both time and space. However, since IGNN is input-invariant, its node representations can be shared across input examples in one batch so that

<sup>1</sup>We assume per-example per-edge per-dimension time cost as a unit time.

$B$  can be removed to get  $O(DT(jEj + jVj))$ . If we sample a smaller set  $\hat{E}$  from  $E$  to run such that  $j\hat{E}j \ll jVj$ , we can further reduce the complexity to  $O(DTjVj)$ .

The pruned message passing in AGNN can be written as

$$\mathbf{H}^t(x) = f_{\text{AGNN}}(\mathbf{H}^{t-1}(x), \mathcal{H}^t; \theta_{G(x)}). \quad (2)$$

Its complexity can be computed similarly as above. However, we cannot remove  $B$ . Fortunately, subgraph  $G(x)$  is not  $G$ . If we let  $x$  be a node  $v$ ,  $G(x)$  grows from a single node, i.e.,  $G^0(x) = \bar{v}v$ , and expands itself each step, leading to a sequence of  $(G^0(x), G^1(x), \dots, G^T(x))$ . Here, we describe the expansion behavior as *consecutive expansion*, which means no jumping across neighborhood allowed, so that we can ensure that

$$G^t(x) \supseteq G^{t-1}(x) \supseteq \partial G^{t-1}(x) \supseteq G^{t-2}(x) \supseteq \partial^2 G^{t-2}(x). \quad (3)$$

Many real-world graphs follow the *small-world* pattern, and the *six degrees of separation* implies  $G^0(x) \supseteq \partial^6 G^0(x) \supseteq G$ . The upper bound of  $G^t(x)$  can grow exponentially in  $t$ , and there is no guarantee that  $G^t(x)$  will not explode.

**Proposition.** *Given a graph  $G$  (undirected or directed in both directions), we assume the probability of the degree of an arbitrary node being less than or equal to  $d$  is larger than  $p$ , i.e.,  $P(\deg(v) \leq d) > p, \forall v \in V$ . Considering a sequence of consecutively expanding subgraphs  $(G^0, G^1, \dots, G^T)$ , starting with  $G^0 = \bar{v}v$ , for all  $t \geq 1$ , we can ensure*

$$P(jV_{G^t}j \geq \frac{d(d-1)^t - 2}{d-2}) > p^{\frac{d(d-1)^t - 2}{d-2}}. \quad (4)$$

The proposition implies the guarantee of upper-bounding  $jV_{G^t(x)}j$  becomes exponentially looser and weaker as  $t$  gets larger even if the given assumption has a small  $d$  and a large  $p$  (close to 1). We define graph increment at step  $t$  as  $\Delta G^t(x)$  such that  $G^t(x) = G^{t-1}(x) \cup \Delta G^t(x)$ . To prevent  $G^t(x)$  from explosion, we need to constrain  $\Delta G^t(x)$ . We propose several sampling strategies:

1.  $\Delta G^t(x) = \hat{\Delta} G^{t-1}(x)$ , which means we sample nodes from the boundary of  $G^{t-1}(x)$ .
2.  $\Delta G^t(x) = \partial \widehat{G^{t-1}(x)}$ , which means we take the boundary of sampled nodes from  $G^{t-1}(x)$ .
3.  $\Delta G^t(x) = \hat{\Delta} \widehat{G^{t-1}(x)}$ , which means we sample nodes from the boundary of  $\widehat{G^{t-1}(x)}$ .
4.  $\Delta G^t(x) = \widehat{\partial \widehat{G^{t-1}(x)}}$ , which means we sample nodes from  $\widehat{\partial \widehat{G^{t-1}(x)}}$ .

Obviously, we have  $\widehat{\partial \widehat{G^{t-1}(x)}} \supseteq \widehat{\Delta} \widehat{G^{t-1}(x)} \supseteq \partial \widehat{G^{t-1}(x)} \supseteq \Delta G^{t-1}(x)$  and  $G^{t-1}(x) \supseteq \partial G^{t-1}(x) \supseteq G^{t-1}(x) \supseteq \partial G^{t-1}(x)$ . Further, we let  $N_g$  and  $N_s$  be the maximum number of sampled nodes in  $\partial \widehat{G^{t-1}(x)}$  and the last sampling of  $\widehat{\Delta} \widehat{G^{t-1}(x)}$  respectively and let  $N_b$  be per-node maximum sampled neighbors in  $\widehat{\Delta} \widehat{G^{t-1}(x)}$ , and then we can obtain much tighter guarantee as follow:

1.  $P(jV_{G^t(x)}j \leq N_g(d-1)) > p^{N_g}$  for  $\partial \widehat{G^{t-1}(x)}$ .
2.  $P(jV_{G^t(x)}j \leq N_g N_b) = 1$  and  $P(jV_{G^t(x)}j \leq N_g \min(d-1, N_b)) > p^{N_g}$  for  $\widehat{\Delta} \widehat{G^{t-1}(x)}$ .
3.  $P(jV_{G^t(x)}j \leq \min(N_g N_b, N_s)) = 1$  for  $\widehat{\Delta} \widehat{G^{t-1}(x)}$ .

By  $\widehat{\Delta} \widehat{G^{t-1}(x)}$ , we can guarantee  $jV_{G^t(x)}j \leq 1 + T \min(N_g N_b, N_s)$ . To constrain the growth of  $G^{t-1}(x)$ , we can decrease either  $N_g N_b$  or  $N_s$ . However, smaller sample size means less area explored and less chance to hit target nodes. We thus use attention operations to do the top- $K$  selection

instead of random sampling when  $K$  has to be small. We change  $\widehat{\Delta} \widehat{G^{t-1}(x)}$  to  $\widehat{\Delta} \widehat{G^{t-1}(x)}$  where  $\widehat{\Delta}$  represents the operation of attending over nodes and picking the top- $K$ . There are two types of attention operations, one applied to  $G^{t-1}(x)$  and the other applied to  $\widehat{\Delta} \widehat{G^{t-1}(x)}$ . Note that the size of  $\widehat{\Delta} \widehat{G^{t-1}(x)}$  might be much larger than  $G^{t-1}(x)$  if we want to sample more nodes with larger  $N_b$  to sufficiently explore the boundary,  $\partial \widehat{G^{t-1}(x)}$ . Nevertheless, we can address this problem by using small dimensions to compute attention scores, since attention carried by each node is just a scalar, much smaller than a node representation vector computed during message passing over  $G^{t-1}(x)$ .

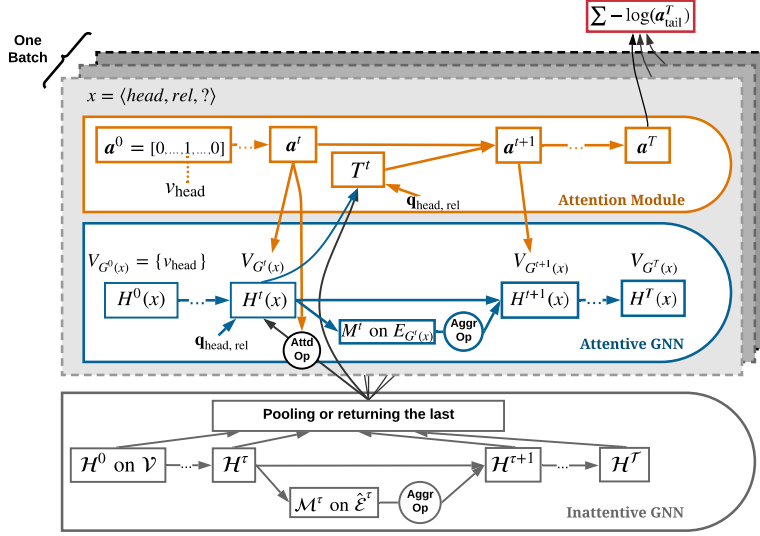


Figure 2: Model architecture used in knowledge graph reasoning.

### 3 MODEL IMPLEMENTATION

#### 3.1 ARCHITECTURE DESIGN FOR KNOWLEDGE GRAPH REASONING

Our model architecture as shown in Figure 2 consists of:

*IGNN module*: performs standard message passing to compute full-graph node representations.

*AGNN module*: performs a batch of pruned message passing to compute input-dependent node representations which also make use of low-level representations from IGNN.

*Attention Module*: performs a flow-style attention transition process, conditioned on node representations from both IGNN and AGNN but only affecting AGNN.

We let  $hV, Ei$  denote a knowledge graph where  $V$  is a set of entities and  $E$  is a set of relations. Each edge or relation is represented by a triple  $hhead, rel, taili$ , where  $head$  is the head entity,  $tail$  is the tail entity, and  $rel$  is their relation type. The goal is to predict potential unknown links, i.e., which entity is likely to be the tail given a query  $hhead, rel, ?i$  with the head and the relation type specified.

**IGNN module.** We implement it using standard message passing mechanism (Gilmer et al., 2017). If the full graph has an extremely large number of edges, we sample a subset of edges,  $\hat{E}^\tau \subseteq E$ , randomly each step. For a batch of input queries, we let node representations from IGNN be shared across queries, containing no batch dimension. Thus, its complexity does not scale with batch size and the saved resources can be allocated to sampling more edges. Each node  $v$  has a state  $\mathcal{H}_{v,:}^\tau$ : at step  $\tau$ , where the initial  $\mathcal{H}_{v,:}^0 = e_v$ . Each edge  $hw^0, r, vi$  produces a message, denoted by  $\mathcal{M}_{hw^0, r, vi,:}^\tau$ : at step  $\tau$ . The computation components include:

Message function:  $\mathcal{M}_{hw^0, r, vi,:}^\tau = \psi_{\text{IGNN}}(\mathcal{H}_{v^0,:}^\tau, e_r, \mathcal{H}_{v,:}^\tau)$ , where  $hw^0, r, vi \in \hat{E}^\tau$ .

Message aggregation:  $\bar{\mathcal{M}}_{v,:}^\tau = \frac{1}{N^\tau(v)} \sum_{v^0, r} \mathcal{M}_{hw^0, r, vi,:}^\tau$ , where  $hw^0, r, vi \in \hat{E}^\tau$ .

Node state update function:  $\mathcal{H}_{v,:}^{\tau+1} = \mathcal{H}_{v,:}^\tau + \delta_{\text{IGNN}}(\mathcal{H}_{v,:}^\tau, \bar{\mathcal{M}}_{v,:}^\tau; e_v)$ , where  $v \in V$ .

We compute messages only for the sampled edges,  $hw^0, r, vi \in \hat{E}^\tau$ , each step. Functions  $\psi_{\text{IGNN}}$  and  $\delta_{\text{IGNN}}$  are implemented by a two-layer MLP (using leakyReLU for the first layer and tanh for the second) with input arguments concatenated respectively. Messages are aggregated by dividing the sum by the square root of  $N^\tau(v)$ , the number of sampled neighbors that send messages to  $v$ , preserving the scale of variance. We use a residual adding to update each node state instead of a GRU or a LSTM. After running for  $T$  steps, we output a pooling result or simply the last, denoted by  $\mathcal{H} = \mathcal{H}^T$ , to feed into downstream modules.

**AGNN module.** AGNN is input-dependent, which means node states depend on input query  $x = \text{head}, \text{rel}, ?i$ , denoted by  $\mathbf{H}_{v,:}^t(x)$ . We implement pruned message passing, running on small subgraphs each conditioned on an input query. We leverage the sparsity and only save  $\mathbf{H}_{v,:}^t(x)$  for visited nodes  $v \in V_{G^t(x)}$ . When  $t = 0$ , we start from node *head* with  $V_{G^0(x)} = \text{head}$ . When computing messages, denoted by  $\mathbf{M}_{h^0,r,vi}^t(x)$ , we use a sampling-attending procedure, explained in Section 3.2, to constrain the number of computed edges. The computation components include:

Message function:  $\mathbf{M}_{h^0,r,vi}^t(x) = \psi_{\text{AGNN}}(\mathbf{H}_{v^0,:}^t(x), \mathbf{c}_r(x), \mathbf{H}_{v,:}^t(x))$ , where  $h^0, r, vi \in E_{G^t(x)}$ , and  $\mathbf{c}_r(x) = [\mathbf{e}_r, \mathbf{q}_{\text{head}}, \mathbf{q}_{\text{rel}}]$  represents a context vector.

Message aggregation:  $\overline{\mathbf{M}}_{v,:}^t(x) = \frac{1}{N^t(v)} \sum_{v^0, r} \mathbf{M}_{h^0,r,vi}^t(x)$ , where  $h^0, r, vi \in E_{G^t(x)}$ .

Node state attending function:  $\widetilde{\mathbf{H}}_{v,:}^{t+1}(x) = a_v^{t+1} \mathbf{W} \mathcal{H}_{v,:}$ , where  $a_v^{t+1}$  is an attention score.

Node state update function:  $\mathbf{H}_{v,:}^{t+1}(x) = \mathbf{H}_{v,:}^t(x) + \delta_{\text{AGNN}}(\mathbf{H}_{v,:}^t(x), \overline{\mathbf{M}}_{v,:}^t(x), \mathbf{c}_v^{t+1}(x))$ , where  $\mathbf{c}_v^{t+1}(x) = [\widetilde{\mathbf{H}}_{v,:}^{t+1}(x), \mathbf{q}_{\text{head}}, \mathbf{q}_{\text{rel}}]$  also represents a context vector.

Query context is defined by its head and relation type embeddings, i.e.,  $\mathbf{q}_{\text{head}} = \mathbf{e}_{\text{head}}$  and  $\mathbf{q}_{\text{rel}} = \mathbf{e}_{\text{rel}}$ . We introduce a node state attending function to pass node representation information from IGNN to AGNN weighted by a scalar attention score  $a_v^{t+1}$  and projected by a learnable matrix  $\mathbf{W}$ . We initialize  $\mathbf{H}_{v,:}^0(x) = \mathcal{H}_{v,:}$  for node  $v \in V_{G^0(x)}$ , treating the rest as zero states.

**Attention module.** Attention over  $T$  steps is represented by a sequence of node probability distributions, denoted by  $\mathbf{a}^t$  ( $t = 1, 2, \dots, T$ ). The initial distribution  $\mathbf{a}^0$  is a one-hot vector with  $\mathbf{a}^0[v_{\text{head}}] = 1$ . To spread attention, we need to compute transition matrices  $\mathbf{T}^t$  each step. Since it is conditioned on both IGNN and AGNN, we capture two types of interaction between  $v^0$  and  $v$ :  $\mathbf{H}_{v^0,:}^t(x) \rightarrow \mathbf{H}_{v,:}^t(x)$ , and  $\mathbf{H}_{v^0,:}^t(x) \rightarrow \mathcal{H}_{v,:}$ . The former favors visited nodes, while the latter is used to attend to unseen nodes.

$$\begin{aligned} \mathbf{T}_{:,v^0}^t &= \text{softmax}_{v \in N^t(v^0)} \left( \sum_r \alpha_1(\mathbf{H}_{v^0,:}^t(x), \mathbf{c}_r(x), \mathbf{H}_{v,:}^t(x)) + \alpha_2(\mathbf{H}_{v^0,:}^t(x), \mathbf{c}_r(x), \mathcal{H}_{v,:}) \right) \\ \alpha_1(\cdot) &= \text{MLP}(\mathbf{H}_{v^0,:}^t(x), \mathbf{c}_r(x))^\top \mathbf{W}_1 \text{MLP}(\mathbf{H}_{v,:}^t(x), \mathbf{c}_r(x)) \\ \alpha_2(\cdot) &= \text{MLP}(\mathbf{H}_{v^0,:}^t(x), \mathbf{c}_r(x))^\top \mathbf{W}_2 \text{MLP}(\mathcal{H}_{v,:}, \mathbf{c}_r(x)) \end{aligned} \quad (5)$$

where  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are two learnable matrices. Each MLP uses one single layer with the leakyReLU activation. To reduce the complexity for computing  $\mathbf{T}^t$ , we use nodes  $v^0 \in \widetilde{V_{G^t(x)}}$ , which contains

nodes with the  $k$ -largest attention scores at step  $t$ , and use nodes  $v$  sampled from  $v^0$ 's neighbors to compute attention transition for the next step. Due to the fact that nodes  $v^0$  result from the top- $k$  pruning, the loss of attention may occur to diminish the total amount. Therefore, we use a renormalized version,  $\mathbf{a}^{t+1} = \mathbf{T}^t \mathbf{a}^t / k \mathbf{T}^t \mathbf{a}^t k$ , to compute new attention scores. We use attention scores at the final step as the probability to predict the tail node.

### 3.2 COMPLEXITY REDUCTION BY ITERATIVE SAMPLING AND ATTENDING

AGNN deals with local subgraphs for each input  $x$  so that only a few selected nodes are kept in  $V_{G^t(x)}$ , called *visited nodes*, and  $|V_{G^t(x)}|$  is much smaller than  $|V|$ . The initial  $V_{G^0(x)}$  contains only one node *head*, and then  $V_{G^t(x)}$  is enlarged each step by adding new nodes. When propagating messages, we can just consider the one-step neighborhood each step. However, the expansion goes so rapidly that it covers almost all nodes after a few steps. The key to address the problem is to constrain the scope of nodes we can expand the boundary from, i.e., the core nodes which determine where we can go next. We call it the *attending-from horizon*,  $\widetilde{G^t(x)}$ , selected according to attention scores  $\mathbf{a}^t$ . Given this horizon, we still need edge sampling over its neighborhood instead of using the whole  $N(\widetilde{G^t(x)})$  in case of a hub node of extremely high degree. Here, we face a trade-off between coverage and complexity when sampling over the neighborhood. Also, we need node representations within each subgraph to keep their information coherent and avoid possible noises caused by randomly sampling. Therefore, we introduce an *attending-to horizon* inside the *sampling*

<sup>2</sup>In practice, we can use a smaller set of edges than  $E_{G^t(x)}$  to pass messages as discussed in Section 3.2

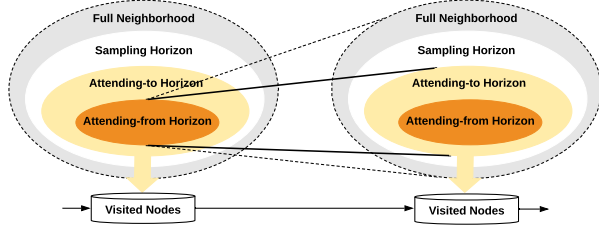


Figure 3: Iterative sampling-attending procedure balancing between coverage and complexity.

*horizon*. We denote the sampling horizon by  $\widehat{N}(G^t(x))$  and the attending-to horizon by  $\widetilde{N}(G^t(x))$ . The attention module runs within the sampling horizon with smaller dimensions in order to sample more neighbors for a larger coverage. Then, we prune the sampling horizon to obtain the attending-to horizon, which contains a subset of nodes selected according to newly computed attention scores  $\mathbf{a}^{t+1}$ . Current message passing iteration at step  $t$  in AGNN can be further constrained on edges between  $G^t(x)$  and  $\widetilde{N}(G^t(x))$ , a smaller set than  $E_{G^t(x)}$ . We illustrate this procedure in Figure 3.

## 4 EXPERIMENTS

**Datasets.** We use six large KG datasets: FB15K, FB15K-237, WN18, WN18RR, NELL995, and YAGO3-10. FB15K-237 (Toutanova & Chen, 2015) is sampled from FB15K (Bordes et al., 2013) with redundant relations removed, and WN18RR (Dettmers et al., 2018) is a subset of WN18 (Bordes et al., 2013) removing triples that cause test leakage. Thus, they are both considered more challenging. NELL995 (Xiong et al., 2017) has separate datasets for 12 query relations each corresponding to a single-query-relation KBC task. YAGO3-10 (Mahdisoltani et al., 2014) contains the largest KG with millions of edges. Their statistics are shown in Table 1. We find some statistical differences between train and validation (or test). In a KG with all training triples as its edges, a triple  $(head, rel, tail)$  is considered as a multi-edge triple if the KG contains other triples that also connect  $head$  and  $tail$  ignoring the direction. We notice that FB15K-237 is a special case compared to the others, as there are no edges in its KG directly linking any pair of  $head$  and  $tail$  in validation (or test). Therefore, when using training triples as queries to train our model, given a batch, for FB15K-237, we cut off from the KG all triples connecting the head-tail pairs in the given batch, ignoring relation types and edge directions, forcing the model to learn a composite reasoning pattern rather than a single-hop pattern, and for the rest datasets, we only remove the triples of this batch and their inverse from the KG to avoid information leakage before training on this batch. This can be regarded as a hyperparameter tuning whether to force a multi-hop reasoning or not, leading to a performance boost of about 2% in HITS@1 on FB15-237.

**Experimental settings.** We use the same data split protocol as in many papers (Dettmers et al., 2018; Xiong et al., 2017; Das et al., 2018). We create a KG, a directed graph, consisting of all train triples and their inverse added for each dataset except NELL995, since it already includes reciprocal relations. Besides, every node in KGs has a self-loop edge to itself. We also add inverse relations into the validation and test set to evaluate the two directions. For evaluation metrics, we use HITS@1,3,10 and the mean reciprocal rank (MRR) in the filtered setting for FB15K-237, WN18RR, FB15K, WN18, and YAGO3-10, and use the mean average precision (MAP) for NELL995’s single-query-relation KBC tasks. For NELL995, we follow the same evaluation procedure as in (Xiong et al., 2017; Das et al., 2018; Shen et al., 2018), ranking the answer entities against the negative examples given in their experiments. We run our experiments using a 12G-memory GPU, TITAN X (Pascal), with Intel(R) Xeon(R) CPU E5-2670 v3 @ 2.30GHz. Our code is written in Python based on TensorFlow 2.0 and NumPy 1.16 and can be found by the link<sup>3</sup> below. We run three times for each hyperparameter setting per dataset to report the means and standard deviations. See hyperparameter details in the appendix.

**Baselines.** We compare our model against embedding-based approaches, including TransE (Bordes et al., 2013), TransR (Lin et al., 2015b), DistMult (Yang et al., 2015), ConvE (Dettmers et al., 2018), ComplE (Trouillon et al., 2016), HolE (Nickel et al., 2016), RotatE (Sun et al., 2018), and ComplEx-N3 (Lacroix et al., 2018), and path-based approaches that use RL methods, including

<sup>3</sup><https://github.com/anonymousauthor123/DPMPN>

Table 1: Statistics of the six KG datasets. PME (tr) means the proportion of multi-edge triples in train; PME (va) means the proportion of multi-edge triples in validation; AL (va) means the average length of shortest paths connecting each head-tail pair in validation.

Dataset	#Entities	#Rels	#Train	#Valid	#Test	PME (tr)	PME (va)	AL (va)
FB15K	14,951	1,345	483,142	50,000	59,071	81.2%	80.6%	1.22
FB15K-237	14,541	237	272,115	17,535	20,466	38.0%	<b>0%</b>	2.25
WN18	40,943	18	141,442	5,000	5,000	93.1%	94.0%	1.18
WN18RR	40,943	11	86,835	3,034	3,134	34.5%	35.5%	2.84
NELL995	74,536	200	149,678	543	2,818	100%	31.1%	2.00
YAGO3-10	123,188	37	1,079,040	5,000	5,000	56.4%	56.0%	1.75

Table 2: Comparison results on the FB15K-237 and WN18RR datasets. Results of [•] are taken from (Nguyen et al., 2018), [/] from (Dettmers et al., 2018), [-] from (Shen et al., 2018), [∩] from (Sun et al., 2018), [∪] from (Das et al., 2018), and [z] from (Lacroix et al., 2018). Some collected results only have a metric score while some including ours take the form of “mean (std)”.

Metric (%)	FB15K-237				WN18RR			
	H@1	H@3	H@10	MRR	H@1	H@3	H@10	MRR
TransE [•]	-	-	46.5	29.4	-	-	50.1	22.6
DistMult [/]	15.5	26.3	41.9	24.1	39	44	49	43
DistMult [-]	20.6 (.4)	31.8 (.2)	-	29.0 (.2)	38.4 (.4)	42.4 (.3)	-	41.3 (.3)
ComplEx [/]	15.8	27.5	42.8	24.7	41	46	51	44
ComplEx [-]	20.8 (.2)	32.6 (.5)	-	29.6 (.2)	38.5 (.3)	43.9 (.3)	-	42.2 (.2)
ConvE [/]	23.7	35.6	50.1	32.5	40	44	52	43
ConvE [-]	23.3 (.4)	33.8 (.3)	-	30.8 (.2)	39.6 (.3)	44.7 (.2)	-	43.3 (.2)
RotatE [∩]	24.1	37.5	53.3	33.8	42.8	49.2	<b>57.1</b>	47.6
ComplEx-N3[z]	-	-	<b>56</b>	<b>37</b>	-	-	57	48
NeuralLP [-]	18.2 (.6)	27.2 (.3)	-	24.9 (.2)	37.2 (.1)	43.4 (.1)	-	43.5 (.1)
MINERVA [-]	14.1 (.2)	23.2 (.4)	-	20.5 (.3)	35.1 (.1)	44.5 (.4)	-	40.9 (.1)
MINERVA [∪]	-	-	45.6	-	41.3	45.6	51.3	-
M-Walk [-]	16.5 (.3)	24.3 (.2)	-	23.2 (.2)	41.4 (.1)	44.5 (.2)	-	43.7 (.1)
<b>DPMPN</b>	<b>28.6 (.1)</b>	<b>40.3 (.1)</b>	53.0 (.3)	36.9 (.1)	<b>44.4 (.4)</b>	<b>49.7 (.8)</b>	55.8 (.5)	<b>48.2 (.5)</b>

DeepPath (Xiong et al., 2017), MINERVA (Das et al., 2018), and M-Walk (Shen et al., 2018), and also that uses learned neural logic, NeuralLP (Yang et al., 2017).

**Comparison results and analysis.** We report comparison on FB15K-23 and WN18RR in Table 2. Our model DPMPN significantly outperforms all the baselines in HITS@1,3 and MRR. Compared to the best baseline, we only lose a few points in HITS@10 but gain a lot in HITS@1,3. We speculate that it is the reasoning capability that helps DPMPN make a sharp prediction by exploiting graph-structured composition locally and conditionally. When a target becomes too vague to predict, reasoning may lose its advantage against embedding-based models. However, path-based baselines, with a certain ability to do reasoning, perform worse than we expect. We argue that it might be inappropriate to think of reasoning, a sequential decision process, equivalent to a sequence of nodes. The average lengths of the shortest paths between heads and tails as shown in Table 1 suggests a very short path, which makes the motivation of using a path almost useless. The reasoning pattern should be modeled in the form of dynamical local graph-structured pattern with nodes densely connected with each other to produce a decision collectively. We also run our model on FB15K, WN18, and YAGO3-10, and the comparison results in the appendix show that DPMPN achieves a very competitive position against the best state of the art. We summarize the comparison on NELL995’s tasks in the appendix. DPMPN performs the best on five tasks, also being competitive on the rest.

**Convergence analysis.** Our model converges very fast during training. We may use half of training queries to train model to generalize as shown in Figure 4(A). Compared to less expensive embedding-based models, our model need to traverse a number of edges when training on one input, consuming much time per batch, but it does not need to pass a second epoch, thus saving a lot of training time. The reason may be that training queries also belong to the KG’s edges and some might be exploited to construct subgraphs during training on other queries.

**Component analysis.** If we do not run message passing in IGNN,  $\mathcal{H}_{v,}$  is just the initial embedding of node  $v$ , and we can still run pruned message passing in AGNN as usual. We want to know whether IGNN is actually useful. Considering that long-range propagated messages might bring in



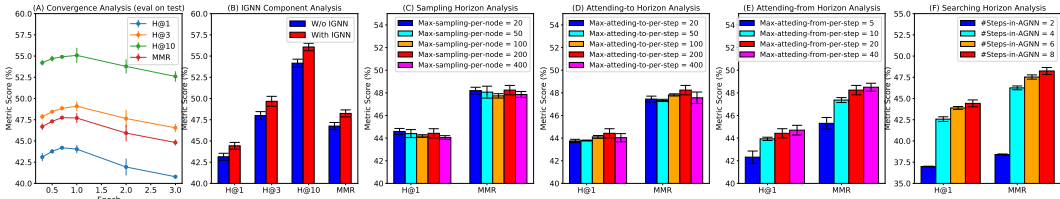


Figure 4: Experimental analysis on WN18RR. (A) Convergence analysis: we pick six model snapshots during training and evaluate them on test. (B) IGNN component analysis: *w/o IGNN* uses zero step to run message passing, while *with IGNN* uses two; (C)-(F) Sampling, attending-to, attending-from and searching horizon analysis. The charts on FB15K-237 can be found in the appendix.

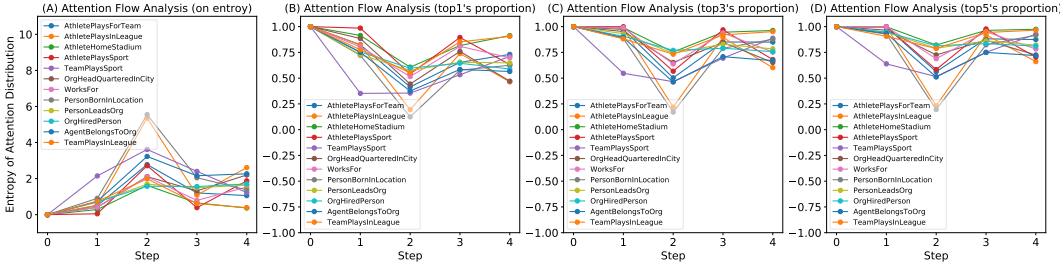


Figure 5: Analysis of attention flow on NELL995 tasks. (A) The average entropy of attention distributions changing along steps for each single-query-relation KBC task. (B)(C)(D) The changing of the proportion of attention concentrated at the top-1,3,5 nodes per step for each task.

noisy features, we compare running IGNN for two steps against totally shutting it down. The result in Figure 4(B) shows that IGNN brings a small gain in each metric on WN18RR.

**Horizon analysis.** The sampling, attending-to, attending-from and searching (i.e., propagation steps) horizons determine how large area a subgraph can expand over. These factors affect computation complexity as well as prediction performance. Intuitively, enlarging the exploring area by sampling more, attending more, and searching longer, may increase the chance of hitting a target to gain some performance. However, the experimental results in Figure 4(C)(D) show that it is not always the case. In Figure 4(E), we can see that increasing the maximum number of attending-from nodes per step is useful, but normal GPUs with a limited memory do not allow for an arbitrarily large number due to heavy intermediate data produced during feedforward computing. Figure 4(F) suggests that the propagation steps of AGNN should not go below four.

**Attention flow analysis.** If the flow-style attention really captures the way we reason about the world, its process should be conducted in a diverging-converging thinking pattern. Intuitively, first, for the diverging thinking phase, we search and collect ideas as much as we can; then, for the converging thinking phase, we try to concentrate our thoughts on one point. To check whether the attention flow has such a pattern, we measure the average entropy of attention distributions changing along steps and also the proportion of attention concentrated at the top-1,3,5 nodes. As we expect, attention is more focused at the final step and the beginning.

**Time cost analysis.** The time cost is affected not only by the scale of a dataset but also by the horizon setting. For each dataset, we list the training time for one epoch corresponding to our standard hyperparameter settings in the appendix. Note that there is always a trade-off between complexity and performance. We thus study whether we can reduce time cost a lot at the price of sacrificing a little performance. We plot the one-epoch training time in Figure 6(A)-(D), using the same settings as we do in the horizon analysis. We can see that *Max-attending-from-per-step* and *#Steps-in-AGNN* affect the training time significantly while *Max-sampling-per-node* and *Max-attending-to-per-step* affect very slightly. Therefore, we can use smaller *Max-sampling-per-node* and *Max-attending-to-per-step* in order to gain a larger batch size, making the computation more efficiency as shown in Figure 6(E).

**Visualization.** To further demonstrate the reasoning capability, we show visualization results of some pruned subgraphs on NELL995’s test data for 12 separate tasks. We avoid using the training data in order to show generalization of the learned reasoning capability. We show the visualization results in Figure 1. See the appendix for detailed analysis and more visualization results.

Figure 6: Analysis of time cost on WN18RR: (A)-(D) measure the one-epoch training time on different horizon settings corresponding to Figure 4(C)-(F); (E) measures on different batch sizes using horizon setting  $\text{Max-sampling-per-node}=20$ ,  $\text{Max-attending-to-per-step}=20$ ,  $\text{Max-attending-from-per-step}=20$ , and  $\text{Steps-in-AGN}=8$ . The charts on FB15K-237 can be found in the appendix.

## 5 RELATED WORK

**Knowledge graph reasoning.** Early work, including TransE (Bordes et al., 2013) and its analogues (Wang et al., 2014; Lin et al., 2015b; Ji et al., 2015), DistMult (Yang et al., 2015), ConvE (Dettmers et al., 2018) and ComplEx (Trouillon et al., 2016), focuses on learning embeddings of entities and relations. Some recent works of this line (Sun et al., 2018; Lacroix et al., 2018) achieve high accuracy. Another line aims to learn inference paths (Lao et al., 2011; Gardner et al., 2014; Guu et al., 2015; Lin et al., 2015a; Toutanova et al., 2016; Das et al., 2017) for knowledge graph reasoning, especially DeepPath (Xiong et al., 2017), MINERVA (Das et al., 2018), and M-Walk (Shen et al., 2018), which use RL to learn multi-hop relational paths. However, these approaches, based on policy gradients or Monte Carlo tree search, often suffer from low sample efficiency and sparse rewards, requiring a large number of rollouts and sophisticated reward function design. Other efforts include learning soft logical rules (Cohen, 2016; Yang et al., 2017) or compositional programs (Liang et al., 2016).

**Relational reasoning in Graph Neural Networks.** Relational reasoning is regarded as the key for combinatorial generalization, taking the form of entity- and relation-centric organization to reason about the composition structure of the world (Craik, 1952; Lake et al., 2017). A multitude of recent implementations (Battaglia et al., 2018) encode relational inductive biases into neural networks to exploit graph-structured representation, including graph convolution networks (GCNs) (Bruna et al., 2014; Henaff et al., 2015; Duvenaud et al., 2015; Kearnes et al., 2016; Defferrard et al., 2016; Niepert et al., 2016; Kipf & Welling, 2017; Bronstein et al., 2017) and graph neural networks (Scarselli et al., 2009; Li et al., 2016; Santoro et al., 2017; Battaglia et al., 2016; Gilmer et al., 2017). Variants of GNN architectures have been developed. Relation networks (Santoro et al., 2017) use a simple but effective neural module to model relational reasoning, and its recurrent versions (Santoro et al., 2018; Palm et al., 2018) do multi-step relational inference for long periods; Interaction networks (Battaglia et al., 2016) provide a general-purpose learnable physics engine, and two of its variants are visual interaction networks (Watters et al., 2017) and vertex attention interaction networks (Hoshen, 2017); Message passing neural networks (Gilmer et al., 2017) unify various GCNs and GNNs into a general message passing formalism by analogy to the one in graphical models.

**Attention mechanism on graphs.** Neighborhood attention operation can enhance GNNs' representation power (Velickovic et al., 2018; Hoshen, 2017; Wang et al., 2018; Kool, 2018). These approaches often use multi-head self-attention to focus on specific interactions with neighbors when aggregating messages, inspired by (Bahdanau et al., 2015; Lin et al., 2017; Vaswani et al., 2017). Most graph-based attention mechanisms attend over neighborhood in a single-hop fashion, and (Hoshen, 2017) claims that the multi-hop architecture does not help to model high-order interaction in experiments. However, a flow-style design of attention in (Xu et al., 2018b) shows a way to model long-range attention, stringing isolated attention operations by transition matrices.

## 6 CONCLUSION

We introduce Dynamically Pruned Message Passing Network (DMPN) and apply it to large-scale knowledge graph reasoning tasks. We propose to learn an input-dependent local subgraph which is progressively and selectively constructed to model a sequential reasoning process in knowledge graphs. We use graphical attention expression, a flow-style attention mechanism, to guide and prune the underlying message passing, making it scalable for large-scale graphs and also providing clear graphical interpretations. We also take the inspiration from the consciousness prior to develop a two-GNN framework to boost experimental performances.

## REFERENCES

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2015.
- Peter W. Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, and Koray Kavukcuoglu. Interaction networks for learning about objects, relations and physics. In *NIPS*, 2016.
- Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, aglar Gülehre, Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.
- Yoshua Bengio. The consciousness prior. *CoRR*, abs/1709.08568, 2017.
- Yoshua Bengio. Challenges for deep learning towards human-level ai, 11 2018.
- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *NIPS*, 2013.
- Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Processing Magazine*, 34:18–42, 2017.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *CoRR*, abs/1312.6203, 2014.
- Wenhu Chen, Wenhan Xiong, Xifeng Yan, and William Yang Wang. Variational knowledge graph reasoning. In *NAACL-HLT*, 2018.
- William W. Cohen. Tensorlog: A differentiable deductive database. *CoRR*, abs/1605.06523, 2016.
- Kenneth H. Craik. The nature of explanation. 1952.
- Rajarshi Das, Arvind Neelakantan, David Belanger, and Andrew McCallum. Chains of reasoning over entities, relations, and text using recurrent neural networks. In *EACL*, 2017.
- Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alexander J. Smola, and Andrew McCallum. Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning. *CoRR*, abs/1711.05851, 2018.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *NIPS*, 2016.
- Stanislas Dehaene, Michel Kerszberg, and Jean Pierre Changeux. A neuronal model of a global workspace in effortful cognitive tasks. *Proceedings of the National Academy of Sciences of the United States of America*, 95 24:14529–34, 1998.
- Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. Convolutional 2d knowledge graph embeddings. In *AAAI*, 2018.
- David K. Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gómez-Bombarelli, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *NIPS*, 2015.
- Matt Gardner, Partha Pratim Talukdar, Jayant Krishnamurthy, and Tom Michael Mitchell. Incorporating vector space similarity in random walk inference over knowledge bases. In *EMNLP*, 2014.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In *ICML*, 2017.

- Kelvin Guu, John Miller, and Percy S. Liang. Traversing knowledge graphs in vector space. In *EMNLP*, 2015.
- Mikael Henaff, Joan Bruna, and Yann LeCun. Deep convolutional networks on graph-structured data. *CoRR*, abs/1506.05163, 2015.
- Yedid Hoshen. Vain: Attentional multi-agent predictive modeling. In *NIPS*, 2017.
- Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jian Zhao. Knowledge graph embedding via dynamic mapping matrix. In *ACL*, 2015.
- Steven M. Kearnes, Kevin McCloskey, Marc Berndl, Vijay S. Pande, and Patrick Riley. Molecular graph convolutions: moving beyond fingerprints. *Journal of computer-aided molecular design*, 30 8:595–608, 2016.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *CoRR*, abs/1609.02907, 2017.
- Wouter Kool. Attention solves your tsp , approximately. 2018.
- Timothée Lacroix, Nicolas Usunier, and Guillaume Obozinski. Canonical tensor decomposition for knowledge base completion. In *ICML*, 2018.
- Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *The Behavioral and brain sciences*, 40:e253, 2017.
- Ni Lao, Tom Michael Mitchell, and William W. Cohen. Random walk inference and learning in a large scale knowledge base. In *EMNLP*, 2011.
- Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. Gated graph sequence neural networks. *CoRR*, abs/1511.05493, 2016.
- Chen Liang, Jonathan Berant, Quoc V. Le, Kenneth D. Forbus, and Ni Lao. Neural symbolic machines: Learning semantic parsers on freebase with weak supervision. In *ACL*, 2016.
- Xi Victoria Lin, Richard Socher, and Caiming Xiong. Multi-hop knowledge graph reasoning with reward shaping. In *EMNLP*, 2018.
- Yankai Lin, Zhiyuan Liu, and Maosong Sun. Modeling relation paths for representation learning of knowledge bases. In *EMNLP*, 2015a.
- Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *AAAI*, 2015b.
- Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *CoRR*, abs/1703.03130, 2017.
- Farzaneh Mahdisoltani, Joanna Asia Biega, and Fabian M. Suchanek. Yago3: A knowledge base from multilingual wikipedias. In *CIDR*, 2014.
- Dai Quoc Nguyen, Tu Dinh Nguyen, Dat Quoc Nguyen, and Dinh Q. Phung. A novel embedding model for knowledge base completion based on convolutional neural network. In *NAACL-HLT*, 2018.
- Maximilian Nickel, Lorenzo Rosasco, and Tomaso A. Poggio. Holographic embeddings of knowledge graphs. In *AAAI*, 2016.
- Mathias Niepert, Mohammed Hassan Ahmed, and Konstantin Kutzkov. Learning convolutional neural networks for graphs. In *ICML*, 2016.
- Rasmus Berg Palm, Ulrich Paquet, and Ole Winther. Recurrent relational networks. In *NeurIPS*, 2018.
- Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic Books, 2018.

- Adam Santoro, David Raposo, David G. T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter W. Battaglia, and Timothy P. Lillicrap. A simple neural network module for relational reasoning. In *NIPS*, 2017.
- Adam Santoro, Ryan Faulkner, David Raposo, Jack W. Rae, Mike Chrzanowski, Théophane Weber, Daan Wierstra, Oriol Vinyals, Razvan Pascanu, and Timothy P. Lillicrap. Relational recurrent neural networks. In *NeurIPS*, 2018.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20:61–80, 2009.
- Yelong Shen, Jianshu Chen, Pu Huang, Yuqing Guo, and Jianfeng Gao. M-walk: Learning to walk over graphs using monte carlo tree search. In *NeurIPS*, 2018.
- Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. *CoRR*, abs/1902.10197, 2018.
- Giulio Tononi, Mélanie Boly, Marcello Massimini, and Christof Koch. Integrated information theory: from consciousness to its physical substrate. *Nature Reviews Neuroscience*, 17:450–461, 2016.
- Kristina Toutanova and Danqi Chen. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, 2015.
- Kristina Toutanova, Victoria Lin, Wen tau Yih, Hoifung Poon, and Chris Quirk. Compositional learning of embeddings for relation paths in knowledge base and text. In *ACL*, 2016.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *ICML*, 2016.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Alejandro Romero, Pietro Lió, and Yoshua Bengio. Graph attention networks. *CoRR*, abs/1710.10903, 2018.
- William Wang. Knowledge graph reasoning: Recent advances, 2018.
- Xiaolong Wang, Ross B. Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7794–7803, 2018.
- Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *AAAI*, 2014.
- Nicholas Watters, Daniel Zoran, Théophane Weber, Peter W. Battaglia, Razvan Pascanu, and Andrea Tacchetti. Visual interaction networks: Learning a physics simulator from video. In *NIPS*, 2017.
- Wenhan Xiong, Thien Hoang, and William Yang Wang. DeepPath: A reinforcement learning method for knowledge graph reasoning. In *EMNLP*, 2017.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *ArXiv*, abs/1810.00826, 2018a.
- Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S. Du, Ken ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *ArXiv*, abs/1905.13211, 2019.
- Xiaoran Xu, Songpeng Zu, Chengliang Gao, Yuan Zhang, and Wei Feng. Modeling attention flow on graphs. *CoRR*, abs/1811.00497, 2018b.
- Bishan Yang, Wen tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. *CoRR*, abs/1412.6575, 2015.
- Fan Yang, Zhilin Yang, and William W. Cohen. Differentiable learning of logical rules for knowledge base reasoning. In *NIPS*, 2017.

## Appendix

### 1 PROOF

**Proposition.** *Given a graph  $G$  (undirected or directed in both directions), we assume the probability of the degree of an arbitrary node being less than or equal to  $d$  is larger than  $p$ , i.e.,  $P(\deg(v) \leq d) > p, \forall v \in V$ . Considering a sequence of consecutively expanding subgraphs  $(G^0, G^1, \dots, G^T)$ , starting with  $G^0 = \{v\}$ , for all  $t \geq 1$ , we can ensure*

$$P(jV_{G^t} \leq \frac{d(d-1)^t - 2}{d-2}) > p^{\frac{d(d-1)^t - 1 - 2}{d-2}}. \quad (6)$$

*Proof.* We consider the extreme case of greedy consecutive expansion, where  $G^t = G^{t-1} \cup \Delta G^t = G^{t-1} \cup \partial G^{t-1}$ , since if this case satisfies the inequality, any case of consecutive expansion can also satisfy it. By definition, all the subgraphs  $G^t$  are a connected graph. Here, we use  $\Delta V^t$  to denote  $V_{G^t}$  for short. In the extreme case, we can ensure that the newly added nodes  $\Delta V^t$  at step  $t$  only belong to the neighborhood of the last added nodes  $\Delta V^{t-1}$ . Since for  $t \geq 2$  each node in  $\Delta V^{t-1}$  already has at least one edge within  $G^{t-1}$  due to the definition of connected graphs, we can have

$$P(j\Delta V^t \leq j\Delta V^{t-1}(d-1)) > p^{j - V^{t-1}j}. \quad (7)$$

For  $t = 1$ , we have  $P(j\Delta V^1 \leq d) > p$  and thus

$$P(jV_{G^1} \leq 1 + d) > p. \quad (8)$$

For  $t \geq 2$ , based on  $jV_{G^t} = 1 + j\Delta V^1 + \dots + j\Delta V^t$ , we obtain

$$P(jV_{G^t} \leq 1 + d + d(d-1) + \dots + d(d-1)^{t-1}) > p^{1+d+d(d-1)+\dots+d(d-1)^{t-2}}, \quad (9)$$

which is

$$P(jV_{G^t} \leq \frac{d(d-1)^t - 2}{d-2}) > p^{\frac{d(d-1)^t - 1 - 2}{d-2}}. \quad (10)$$

We can find that  $t = 1$  also satisfies this inequality.  $\square$

## 2 HYPERPARAMETER SETTINGS

Table 3: Our standard hyperparameter settings we use for each dataset plus their one-epoch training time. For experimental analysis, we only adjust one hyperparameter and keep the remaining fixed as the standard setting. For NELL995, the one-epoch training time means the average time cost of the 12 single-query-relation tasks.

Hyperparameter	FB15K-237	FB15K	WN18RR	WN18	YAGO3-10	NELL995
<i>batch_size</i>	80	80	100	100	100	10
<i>n_dims_att</i>	50	50	50	50	50	200
<i>n_dims</i>	100	100	100	100	100	200
<i>max_sampling_per_step (in IGNN)</i>	10000	10000	10000	10000	10000	10000
<i>max_attending_from_per_step</i>	20	20	20	20	20	100
<i>max_sampling_per_node (in AGNN)</i>	200	200	200	200	200	1000
<i>max_attending_to_per_step</i>	200	200	200	200	200	1000
<i>n_steps_in_IGNN</i>	2	1	2	1	1	1
<i>n_steps_in_AGNN</i>	6	6	8	8	6	5
<i>learning_rate</i>	0.001	0.001	0.001	0.001	0.0001	0.001
<i>optimizer</i>	Adam	Adam	Adam	Adam	Adam	Adam
<i>grad_clipnorm</i>	1	1	1	1	1	1
<i>n_epochs</i>	1	1	1	1	1	3
One-epoch training time (h)	25.7	63.7	4.3	8.5	185.0	0.12

The hyperparameters can be categorized into three groups:

Normal hyperparameters, including *batch\_size*, *n\_dims\_att*, *n\_dims*, *learning\_rate*, *grad\_clipnorm*, and *n\_epochs*. We set smaller dimensions, *n\_dims\_att*, for computation in the attention module, as it uses more edges than the message passing uses in AGNN, and also intuitively, it does not need to propagate high-dimensional messages but only compute scalar scores over a sampled neighborhood, in concert with the idea in the key-value mechanism (Bengio, 2017). We set *n\_epochs* = 1 in most cases, indicating that our model can be trained well by one epoch only due to its fast convergence.

The hyperparameters in charge of the sampling-attending horizon, including *max\_sampling\_per\_step* that controls the maximum number to sample edges per step in IGNN, and *max\_sampling\_per\_node*, *max\_attending\_from\_per\_step* and *max\_attending\_to\_per\_step* that control the maximum number to sample neighbors of each selected node per step per input, the maximum number of selected nodes for attending-from per step per input, and the maximum number of selected nodes in a sampled neighborhood for attending-to per step per input in AGNN.

The hyperparameters in charge of the searching horizon, including *n\_steps\_in\_IGNN* representing the number of propagation steps to run standard message passing in IGNN, and *n\_steps\_in\_AGNN* representing the number of propagation steps to run pruned message passing in AGNN.

Note that we tune these hyperparameters according to not only their performances but also the computation resources available to us. In some cases, to deal with a very large knowledge graph with limited resources, we need to make a trade-off between efficiency and effectiveness. For example, each of NELL995’s single-query-relation tasks has a small training set, though still with a large graph, so we can reduce the batch size in favor of affording larger dimensions and a larger sampling-attending horizon without any concern for waiting too long to finish one epoch.

### 3 MORE EXPERIMENTAL RESULTS

Table 4: Comparison results on the FB15K and WN18 datasets. Results of [•] are taken from (Nickel et al., 2016), [/] from (Dettmers et al., 2018), [j] from (Sun et al., 2018), [-] from (Yang et al., 2017), and [z] from (Lacroix et al., 2018). Our results take the form of "mean (std)".

Metric (%)	FB15K				WN18			
	H@1	H@3	H@10	MRR	H@1	H@3	H@10	MRR
TransE [•]	29.7	57.8	74.9	46.3	11.3	88.8	94.3	49.5
HoIE [•]	40.2	61.3	73.9	52.4	93.0	94.5	94.9	93.8
DistMult [/]	54.6	73.3	82.4	65.4	72.8	91.4	93.6	82.2
ComplEx [/]	59.9	75.9	84.0	69.2	93.6	93.6	94.7	94.1
ConvE [/]	55.8	72.3	83.1	65.7	93.5	94.6	95.6	94.3
RotatE [j]	<b>74.6</b>	<b>83.0</b>	88.4	79.7	<b>94.4</b>	<b>95.2</b>	95.9	94.9
ComplEx-N3 [z]	-	-	<b>91</b>	<b>86</b>	-	-	96	95
NeuralLP [-]	-	-	83.7	76	-	-	94.5	94
<b>DPMPN</b>	72.6 (.4)	78.4 (.4)	83.4 (.5)	76.4 (.4)	91.6 (.8)	93.6 (.4)	94.9 (.4)	92.8 (.6)

Table 5: Comparison results on the YAGO3-10 dataset. Results of [•] are taken from (Dettmers et al., 2018), [/] from (Lacroix et al., 2018), and [z] from (Lacroix et al., 2018).

Metric (%)	YAGO3-10			
	H@1	H@3	H@10	MRR
DistMult [•]	24	38	54	34
ComplEx [•]	26	40	55	36
ConvE [•]	35	49	62	44
ComplEx-N3 [z]	-	-	<b>71</b>	<b>58</b>
<b>DPMPN</b>	<b>48.4</b>	<b>59.5</b>	67.9	55.3

Table 6: Comparison results of MAP scores (%) on NELL995’s single-query-relation KBC tasks. We take our baselines’ results from (Shen et al., 2018). No reports found on the last two in the paper.

Tasks	NeuCFlow	M-Walk	MINERVA	DeepPath	TransE	TransR
AthletePlaysForTeam	83.9 (0.5)	<b>84.7 (1.3)</b>	82.7 (0.8)	72.1 (1.2)	62.7	67.3
AthletePlaysInLeague	97.5 (0.1)	<b>97.8 (0.2)</b>	95.2 (0.8)	92.7 (5.3)	77.3	91.2
AthleteHomeStadium	<b>93.6 (0.1)</b>	91.9 (0.1)	92.8 (0.1)	84.6 (0.8)	71.8	72.2
AthletePlaysSport	<b>98.6 (0.0)</b>	98.3 (0.1)	<b>98.6 (0.1)</b>	91.7 (4.1)	87.6	96.3
TeamPlayssport	<b>90.4 (0.4)</b>	88.4 (1.8)	87.5 (0.5)	69.6 (6.7)	76.1	81.4
OrgHeadQuarteredInCity	94.7 (0.3)	<b>95.0 (0.7)</b>	94.5 (0.3)	79.0 (0.0)	62.0	65.7
WorksFor	<b>86.8 (0.0)</b>	84.2 (0.6)	82.7 (0.5)	69.9 (0.3)	67.7	69.2
PersonBornInLocation	<b>84.1 (0.5)</b>	81.2 (0.0)	78.2 (0.0)	75.5 (0.5)	71.2	81.2
PersonLeadsOrg	88.4 (0.1)	<b>88.8 (0.5)</b>	83.0 (2.6)	79.0 (1.0)	75.1	77.2
OrgHiredPerson	84.7 (0.8)	<b>88.8 (0.6)</b>	87.0 (0.3)	73.8 (1.9)	71.9	73.7
AgentBelongsToOrg	<b>89.3 (1.2)</b>	-	-	-	-	-
TeamPlaysInLeague	<b>97.2 (0.3)</b>	-	-	-	-	-



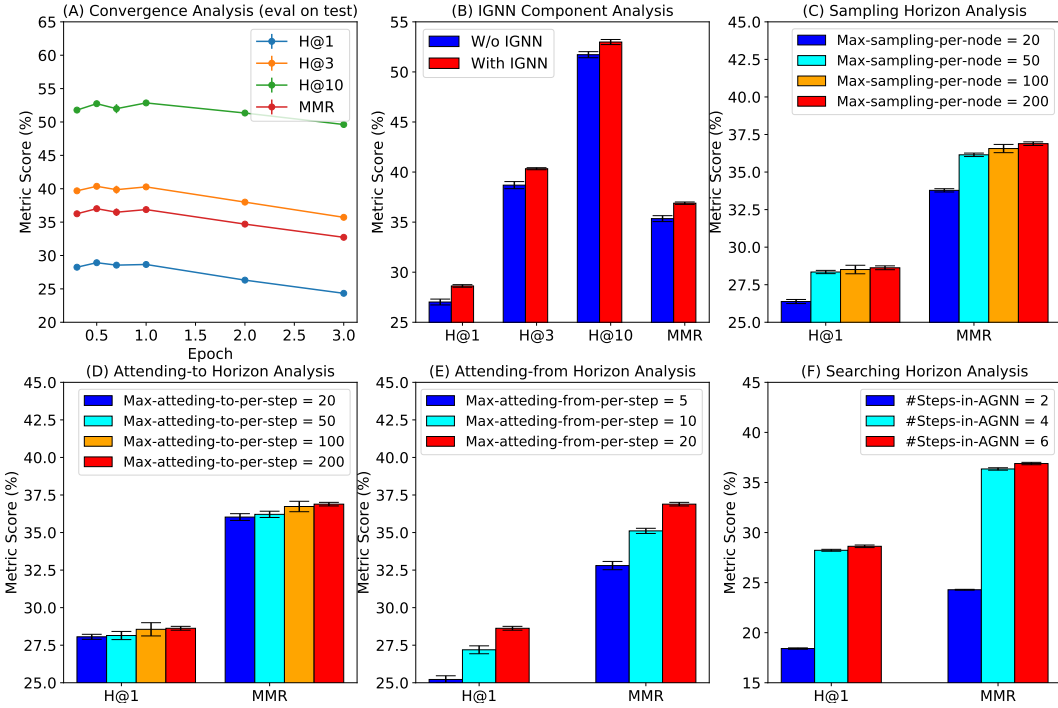


Figure 7: Experimental analysis on FB15K-237. (A) Convergence analysis: we pick six model snapshots at time points of 0.3, 0.5, 0.7, 1, 2, and 3 epochs during training and evaluate them on test; (B) IGNN component analysis: *w/o IGNN* uses zero step to run message passing, while *with IGNN* uses two steps; (C)-(F) Sampling, attending-to, attending-from and searching horizon analysis.

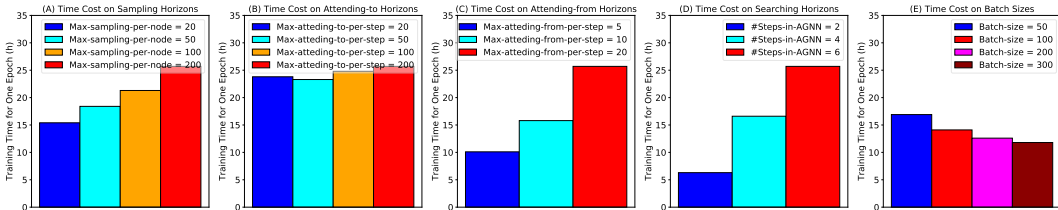


Figure 8: Analysis of time cost on FB15K-237: (A)-(D) measure the one-epoch training time on different horizon settings corresponding to Figure 7(C)-(F); (E) measures on different batch sizes using horizon setting *Max-sampled-edges-per-node=20*, *Max-seen-nodes-per-step=20*, *Max-attended-nodes-per-step=20*, and *#Steps-of-AGNN=6*.

## 4 MORE VISUALIZATION RESULTS

### 4.1 CASE STUDY ON THE ATHLETEPLAYSFORTTEAM TASK

In the case shown in Figure 9, the query is (*concept-personnorthamerica\_michael\_turner*, *concept-athleteplaysforteam*, ?) and a true answer is *concept-sportsteam\_falcons*. From Figure 9, we can see our model learns that (*concept-personnorthamerica\_michael\_turner*, *concept-athlethomestadium*, *concept-stadiumeventvenue\_georgia\_dome*) and (*concept-stadiumeventvenue\_georgia\_dome*, *concept:teahomestadium\_inv*, *concept-sportsteam\_falcons*) are two important facts to support the answer of *concept-sportsteam\_falcons*. Besides, other facts, such as (*concept-athlete\_joey\_harrington*, *concept-athlethomestadium*, *concept-stadiumeventvenue\_georgia\_dome*) and (*concept-athlete\_joey\_harrington*, *concept:athleteplaysforteam*, *concept-sportsteam\_falcons*), provide a vivid example that a person or an athlete with *concept-stadiumeventvenue\_georgia\_dome* as his or her home stadium might play for the team *concept-sportsteam\_falcons*. We have such examples more than one, like *concept-athlete\_rodny\_white*'s and *concept-athlete\_quarterback\_matt\_ryan*'s. The entity *con-*

`cept_sportsleague_nfl` cannot help us differentiate the true answer from other NFL teams, but it can at least exclude those non-NFL teams. In a word, our subgraph-structured representation can well capture the relational and compositional reasoning pattern.

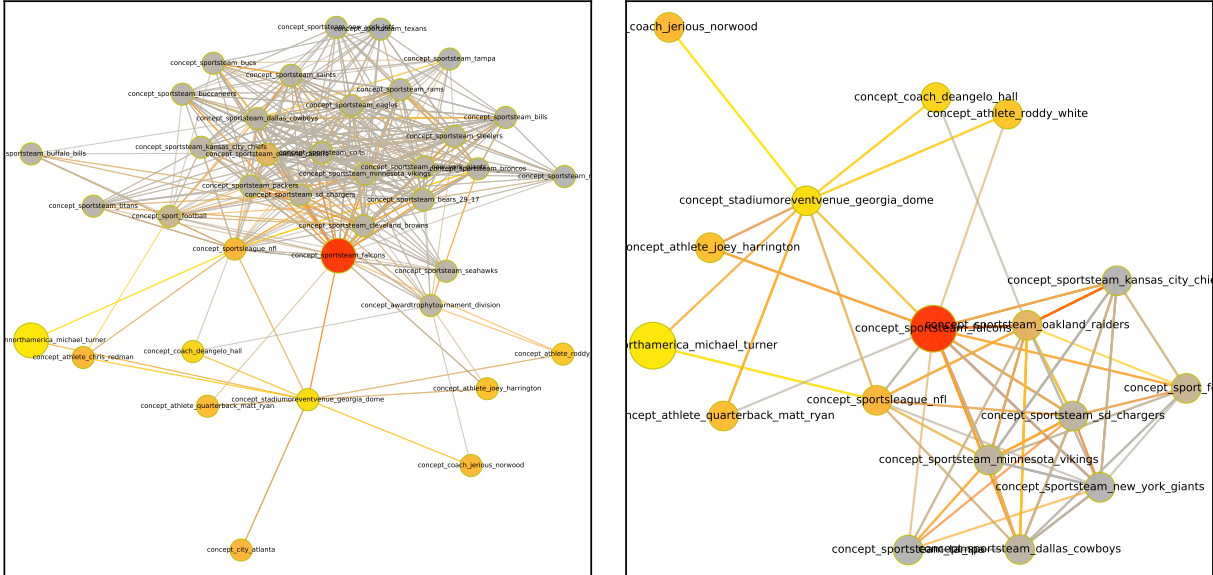


Figure 9: **AthletePlaysForTeam**. The head is `concept-personnorthamerica_michael_turner`, the query relation is `concept:athleteplaysforteam`, and the tail is `concept-sportsteam_falcons`. The left is a full subgraph derived with `max_attending_from_per_step=20`, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the AthletePlaysForTeam task

Query: (`concept-personnorthamerica_michael_turner`, `concept:athleteplaysforteam`, `concept-sportsteam_falcons`)

Selected key edges:

`concept-personnorthamerica_michael_turner`, `concept:agentbelongstoorganization`, `concept_sportsleague_nfl`  
`concept-personnorthamerica_michael_turner`, `concept:athlethomestadium`, `concept_stadiumoreventvenue_georgia_dome`  
`concept_sportsleague_nfl`, `concept:agentcompeteswithagent`, `concept_sportsleague_nfl`  
`concept_sportsleague_nfl`, `concept:agentcompeteswithagent_inv`, `concept_sportsleague_nfl`  
`concept_sportsleague_nfl`, `concept:teamplaysinleague_inv`, `concept_sportsteam_sd_chargers`  
`concept_sportsleague_nfl`, `concept:leaguestadiums`, `concept_stadiumoreventvenue_georgia_dome`  
`concept_sportsleague_nfl`, `concept:teamplaysinleague_inv`, `concept_sportsteam_falcons`  
`concept_sportsleague_nfl`, `concept:agentbelongstoorganization_inv`, `concept-personnorthamerica_michael_turner`  
`concept_stadiumoreventvenue_georgia_dome`, `concept:teamhomestadium_inv`, `concept_sportsteam_falcons`  
`concept_stadiumoreventvenue_georgia_dome`, `concept:athlethomestadium_inv`, `concept_athlete_joey_harrington`  
`concept_stadiumoreventvenue_georgia_dome`, `concept:athlethomestadium_inv`, `concept_athlete_rodny_white`  
`concept_stadiumoreventvenue_georgia_dome`, `concept:athlethomestadium_inv`, `concept_coach_deangelo_hall`  
`concept_stadiumoreventvenue_georgia_dome`, `concept:athlethomestadium_inv`, `concept-personnorthamerica_michael_turner`  
`concept_sportsleague_nfl`, `concept:subpartoforganization_inv`, `concept_sportsteam_oakland_raiders`  
`concept_sportsteam_sd_chargers`, `concept:teamplaysinleague`, `concept_sportsleague_nfl`  
`concept_sportsteam_sd_chargers`, `concept:teamplaysagainstteam`, `concept_sportsteam_falcons`  
`concept_sportsteam_sd_chargers`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_falcons`  
`concept_sportsteam_sd_chargers`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_oakland_raiders`  
`concept_sportsteam_sd_chargers`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_oakland_raiders`  
`concept_sportsteam_falcons`, `concept:teamplaysinleague`, `concept_sportsleague_nfl`  
`concept_sportsteam_falcons`, `concept:teamplaysagainstteam`, `concept_sportsteam_sd_chargers`  
`concept_sportsteam_falcons`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_sd_chargers`  
`concept_sportsteam_falcons`, `concept:teamhomestadium`, `concept_stadiumoreventvenue_georgia_dome`  
`concept_sportsteam_falcons`, `concept:teamplaysagainstteam`, `concept_sportsteam_oakland_raiders`  
`concept_sportsteam_falcons`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_oakland_raiders`  
`concept_sportsteam_falcons`, `concept:athleleledsportsteam_inv`, `concept_athlete_joey_harrington`  
`concept_athlete_joey_harrington`, `concept:athlethomestadium`, `concept_stadiumoreventvenue_georgia_dome`  
`concept_athlete_joey_harrington`, `concept:athleleledsportsteam`, `concept_sportsteam_falcons`  
`concept_athlete_rodny_white`, `concept:athlethomestadium`, `concept_stadiumoreventvenue_georgia_dome`  
`concept_athlete_rodny_white`, `concept:athleteplaysforteam`, `concept_sportsteam_falcons`  
`concept_coach_deangelo_hall`, `concept:athlethomestadium`, `concept_stadiumoreventvenue_georgia_dome`

concept\_coach\_deangelo\_hall , concept\_athleteplaysforteam , concept\_sportsteam\_oakland\_raiders  
 concept\_sportsleague\_nfl , concept\_teamplaysinleague\_inv , concept\_sportsteam\_new\_york\_giants  
 concept\_sportsteam\_sd\_chargers , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_new\_york\_giants  
 concept\_sportsteam\_falcons , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_new\_york\_giants  
 concept\_sportsteam\_oakland\_raiders , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_new\_york\_giants  
 concept\_sportsteam\_oakland\_raiders , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_sd\_chargers  
 concept\_sportsteam\_oakland\_raiders , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_sd\_chargers  
 concept\_sportsteam\_oakland\_raiders , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_falcons  
 concept\_sportsteam\_oakland\_raiders , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_falcons  
 concept\_sportsteam\_oakland\_raiders , concept\_agentcompeteswithagent\_inv , concept\_sportsteam\_oakland\_raiders  
 concept\_sportsteam\_new\_york\_giants , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_sd\_chargers  
 concept\_sportsteam\_new\_york\_giants , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_falcons  
 concept\_sportsteam\_new\_york\_giants , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_falcons  
 concept\_sportsteam\_new\_york\_giants , concept\_teamplaysagainstteam\_inv , concept\_sportsteam\_oakland\_raiders

## 4.2 MORE RESULTS

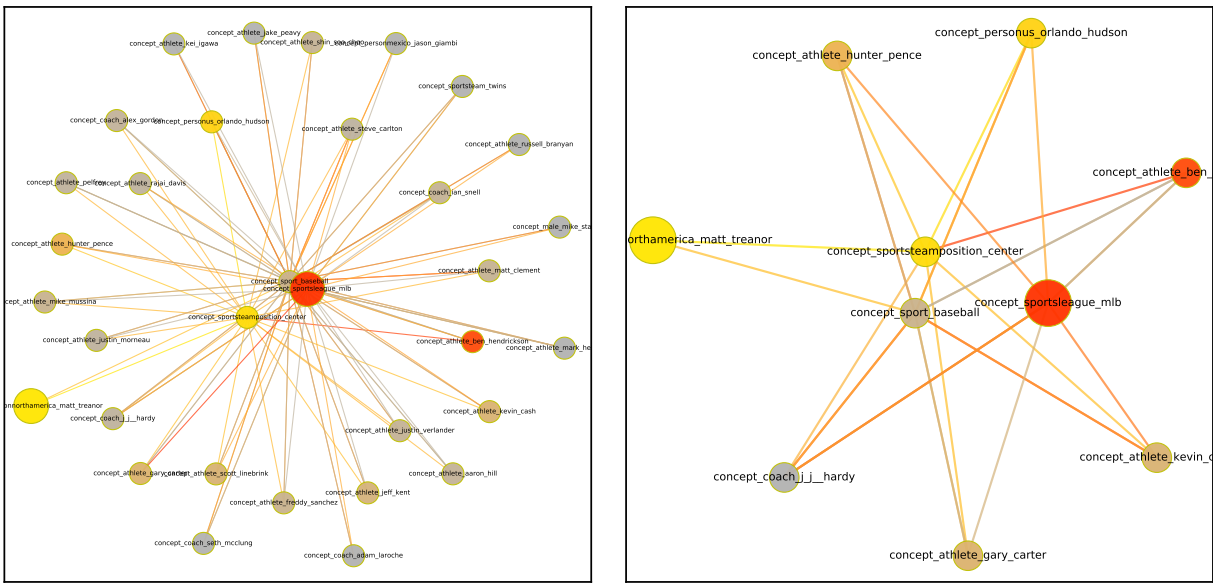


Figure 10: **AthletePlaysInLeague**. The head is *concept\_personnorthamerica\_matt\_treanor*, the query relation is *concept\_athleteplaysinleague*, and the tail is *concept\_sportsleague\_mlb*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the AthletePlaysInLeague task

Query: (*concept\_personnorthamerica\_matt\_treanor* , *concept\_athleteplaysinleague* , *concept\_sportsleague\_mlb*)

Selected key edges:

*concept\_personnorthamerica\_matt\_treanor* , *concept\_athleteflyouttosportsteamposition* , *concept\_sportsteamposition\_center*  
*concept\_personnorthamerica\_matt\_treanor* , *concept\_athleteplayssport* , *concept\_sport\_baseball*  
*concept\_sportsteamposition\_center* , *concept\_athleteflyouttosportsteamposition\_inv* , *concept\_personus\_orlando\_hudson*  
*concept\_sportsteamposition\_center* , *concept\_athleteflyouttosportsteamposition\_inv* , *concept\_athlete\_ben\_hendrickson*  
*concept\_sportsteamposition\_center* , *concept\_athleteflyouttosportsteamposition\_inv* , *concept\_coach\_j\_j\_hardy*  
*concept\_sportsteamposition\_center* , *concept\_athleteflyouttosportsteamposition\_inv* , *concept\_athlete\_hunter\_pence*  
*concept\_sport\_baseball* , *concept\_athleteplayssport\_inv* , *concept\_personus\_orlando\_hudson*  
*concept\_sport\_baseball* , *concept\_athleteplayssport\_inv* , *concept\_athlete\_ben\_hendrickson*  
*concept\_sport\_baseball* , *concept\_athleteplayssport\_inv* , *concept\_coach\_j\_j\_hardy*  
*concept\_sport\_baseball* , *concept\_athleteplayssport\_inv* , *concept\_athlete\_hunter\_pence*  
*concept\_personus\_orlando\_hudson* , *concept\_athleteplaysinleague* , *concept\_sportsleague\_mlb*  
*concept\_personus\_orlando\_hudson* , *concept\_athleteplayssport* , *concept\_sport\_baseball*  
*concept\_athlete\_ben\_hendrickson* , *concept\_coachesinleague* , *concept\_sportsleague\_mlb*  
*concept\_athlete\_ben\_hendrickson* , *concept\_athleteplayssport* , *concept\_sport\_baseball*

concept\_coach\_j\_j\_hardy, concept:coachesinleague, concept:sportsleague\_mlb  
 concept\_coach\_j\_j\_hardy, concept:athleteplaysinleague, concept:sportsleague\_mlb  
 concept\_coach\_j\_j\_hardy, concept:athleteplayssport, concept:sport\_baseball  
 concept\_athlete\_hunter\_pence, concept:athleteplaysinleague, concept:sportsleague\_mlb  
 concept\_athlete\_hunter\_pence, concept:athleteplayssport, concept:sport\_baseball  
 concept\_sportsleague\_mlb, concept:coachesinleague\_inv, concept\_athlete\_ben\_hendrickson  
 concept\_sportsleague\_mlb, concept:coachesinleague\_inv, concept\_coach\_j\_j\_hardy

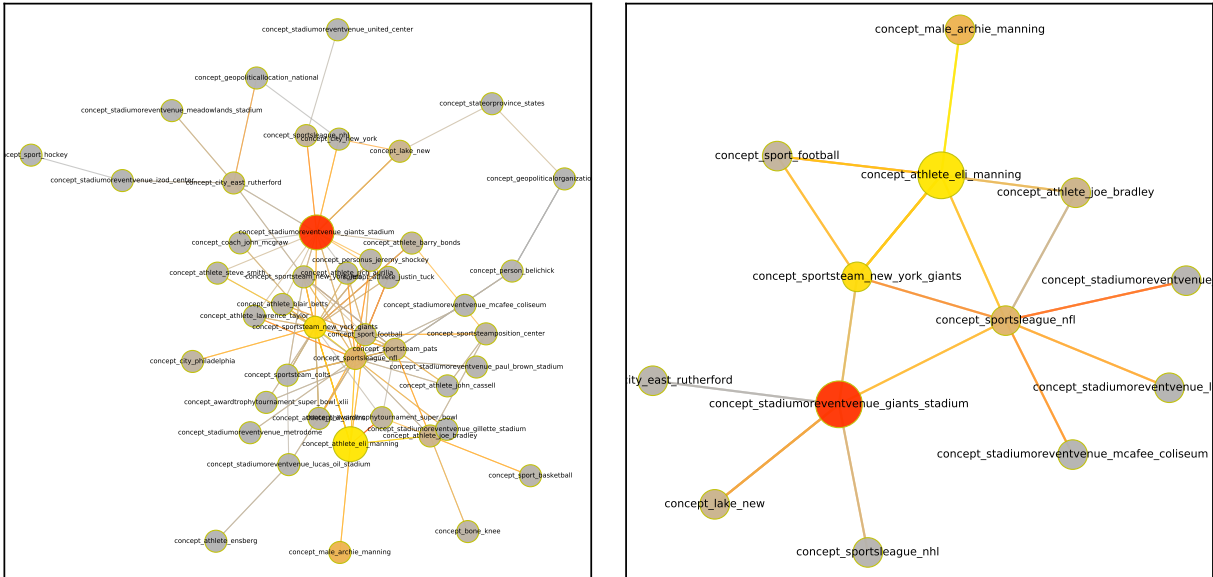


Figure 11: **AthleteHomeStadium**. The head is *concept\_athlete\_eli\_manning*, the query relation is *concept:athletehomestadium*, and the tail is *concept\_stadiumoreventvenue\_giants\_stadium*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the AthleteHomeStadium task

Query: (concept.athlete\_eli\_manning, concept:athletehomestadium, concept\_stadiumoreventvenue\_giants\_stadium)

Selected key edges:

concept\_athlete\_eli\_manning, concept:personbelongstoorganization, concept\_sportsteam\_new\_york\_giants  
 concept\_athlete\_eli\_manning, concept:athleteplaysforteam, concept\_sportsteam\_new\_york\_giants  
 concept\_athlete\_eli\_manning, concept:athleteledsportsteam, concept\_sportsteam\_new\_york\_giants  
 concept\_athlete\_eli\_manning, concept:athleteplaysinleague, concept\_sportsleague\_nfl  
 concept\_athlete\_eli\_manning, concept:fatherofperson\_inv, concept\_male\_archie\_manning  
 concept\_sportsteam\_new\_york\_giants, concept:teamplaysinleague, concept\_sportsleague\_nfl  
 concept\_sportsteam\_new\_york\_giants, concept:teamhomestadium, concept\_stadiumoreventvenue\_giants\_stadium  
 concept\_sportsteam\_new\_york\_giants, concept:personbelongstoorganization\_inv, concept\_athlete\_eli\_manning  
 concept\_sportsteam\_new\_york\_giants, concept:athleteplaysforteam\_inv, concept\_athlete\_eli\_manning  
 concept\_sportsteam\_new\_york\_giants, concept:athleteledsportsteam\_inv, concept\_athlete\_eli\_manning  
 concept\_sportsleague\_nfl, concept:teamplaysinleague\_inv, concept\_sportsteam\_new\_york\_giants  
 concept\_sportsleague\_nfl, concept:agentcompeteswithagent, concept\_sportsleague\_nfl  
 concept\_sportsleague\_nfl, concept:agentcompeteswithagent\_inv, concept\_sportsleague\_nfl  
 concept\_sportsleague\_nfl, concept:leaguestadiums, concept\_stadiumoreventvenue\_giants\_stadium  
 concept\_sportsleague\_nfl, concept:athleteplaysinleague\_inv, concept\_athlete\_eli\_manning  
 concept\_male\_archie\_manning, concept:fatherofperson, concept\_athlete\_eli\_manning  
 concept\_sportsleague\_nfl, concept:leaguestadiums, concept\_stadiumoreventvenue\_paul\_brown\_stadium  
 concept\_stadiumoreventvenue\_giants\_stadium, concept:teamhomestadium\_inv, concept\_sportsteam\_new\_york\_giants  
 concept\_stadiumoreventvenue\_giants\_stadium, concept:leaguestadiums\_inv, concept\_sportsleague\_nfl  
 concept\_stadiumoreventvenue\_giants\_stadium, concept:proxyfor\_inv, concept\_city\_east\_rutherford  
 concept\_city\_east\_rutherford, concept:proxyfor, concept\_stadiumoreventvenue\_giants\_stadium  
 concept\_stadiumoreventvenue\_paul\_brown\_stadium, concept:leaguestadiums\_inv, concept\_sportsleague\_nfl

### For the AthletePlaysSport task

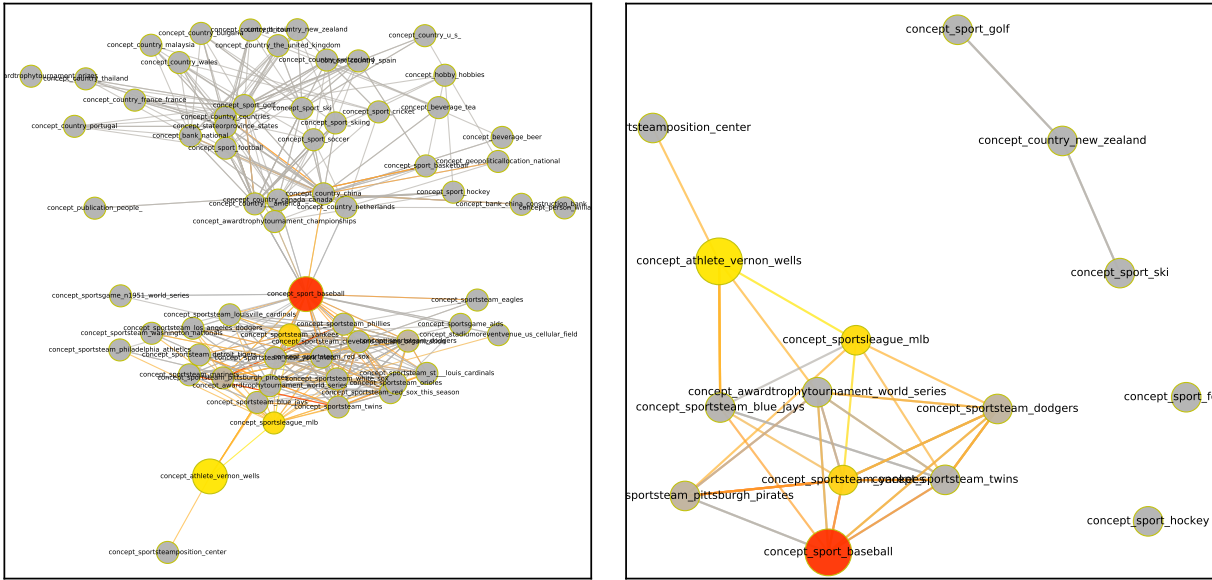


Figure 12: **AthletePlaysSport**. The head is `concept_athlete_vernon_wells`, the query relation is `concept:athleteplaysport`, and the tail is `concept_sport_baseball`. The left is a full subgraph derived with `max_attending_from_per_step=20`, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

Query: (`concept_athlete_vernon_wells`, `concept:athleteplaysport`, `concept_sport_baseball`)

Selected key edges:

`concept_athlete_vernon_wells`, `concept:athleteplaysinleague`, `concept_sportsleague_mlb`  
`concept_athlete_vernon_wells`, `concept:coachwontrophy`, `concept_awardtrophytournament_world_series`  
`concept_athlete_vernon_wells`, `concept:agentcollaborateswithagent_inv`, `concept_sportsteam_blue_jays`  
`concept_athlete_vernon_wells`, `concept:personbelongstoorganization`, `concept_sportsteam_blue_jays`  
`concept_athlete_vernon_wells`, `concept:athleteplaysforteam`, `concept_sportsteam_blue_jays`  
`concept_athlete_vernon_wells`, `concept:athleteledsportsteam`, `concept_sportsteam_blue_jays`  
`concept_sportsleague_mlb`, `concept:teamplaysinleague_inv`, `concept_sportsteam_dodgers`  
`concept_sportsleague_mlb`, `concept:teamplaysinleague_inv`, `concept_sportsteam_yankees`  
`concept_sportsleague_mlb`, `concept:teamplaysinleague_inv`, `concept_sportsteam_pittsburgh_pirates`  
`concept_awardtrophytournament_world_series`, `concept:teamwontrophy_inv`, `concept_sportsteam_dodgers`  
`concept_awardtrophytournament_world_series`, `concept:teamwontrophy_inv`, `concept_sportsteam_yankees`  
`concept_awardtrophytournament_world_series`, `concept:awardtrophytournamentisthechampionshipgameofthenationalsport`,  
`concept_sport_baseball`  
`concept_awardtrophytournament_world_series`, `concept:teamwontrophy_inv`, `concept_sportsteam_pittsburgh_pirates`  
`concept_sportsteam_blue_jays`, `concept:teamplaysinleague`, `concept_sportsleague_mlb`  
`concept_sportsteam_blue_jays`, `concept:teamplaysagainstteam`, `concept_sportsteam_yankees`  
`concept_sportsteam_blue_jays`, `concept:teamplayssport`, `concept_sport_baseball`  
`concept_sportsteam_dodgers`, `concept:teamplaysagainstteam`, `concept_sportsteam_yankees`  
`concept_sportsteam_dodgers`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_yankees`  
`concept_sportsteam_dodgers`, `concept:teamwontrophy`, `concept_awardtrophytournament_world_series`  
`concept_sportsteam_dodgers`, `concept:teamplayssport`, `concept_sport_baseball`  
`concept_sportsteam_yankees`, `concept:teamplaysagainstteam`, `concept_sportsteam_dodgers`  
`concept_sportsteam_yankees`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_dodgers`  
`concept_sportsteam_yankees`, `concept:teamwontrophy`, `concept_awardtrophytournament_world_series`  
`concept_sportsteam_yankees`, `concept:teamplayssport`, `concept_sport_baseball`  
`concept_sportsteam_yankees`, `concept:teamplaysagainstteam`, `concept_sportsteam_pittsburgh_pirates`  
`concept_sportsteam_yankees`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_pittsburgh_pirates`  
`concept_sport_baseball`, `concept:teamplayssport_inv`, `concept_sportsteam_dodgers`  
`concept_sport_baseball`, `concept:teamplayssport_inv`, `concept_sportsteam_yankees`  
`concept_sport_baseball`, `concept:awardtrophytournamentisthechampionshipgameofthenationalsport_inv`,  
`concept_awardtrophytournament_world_series`  
`concept_sport_baseball`, `concept:teamplayssport_inv`, `concept_sportsteam_pittsburgh_pirates`  
`concept_sportsteam_pittsburgh_pirates`, `concept:teamplaysagainstteam`, `concept_sportsteam_yankees`  
`concept_sportsteam_pittsburgh_pirates`, `concept:teamplaysagainstteam_inv`, `concept_sportsteam_yankees`

concept\_sportsteam\_pittsburgh\_pirates , concept:teamwontrophy , concept\_awardtrophytournament\_world\_series  
 concept\_sportsteam\_pittsburgh\_pirates , concept:teamplayssport , concept\_sport\_baseball

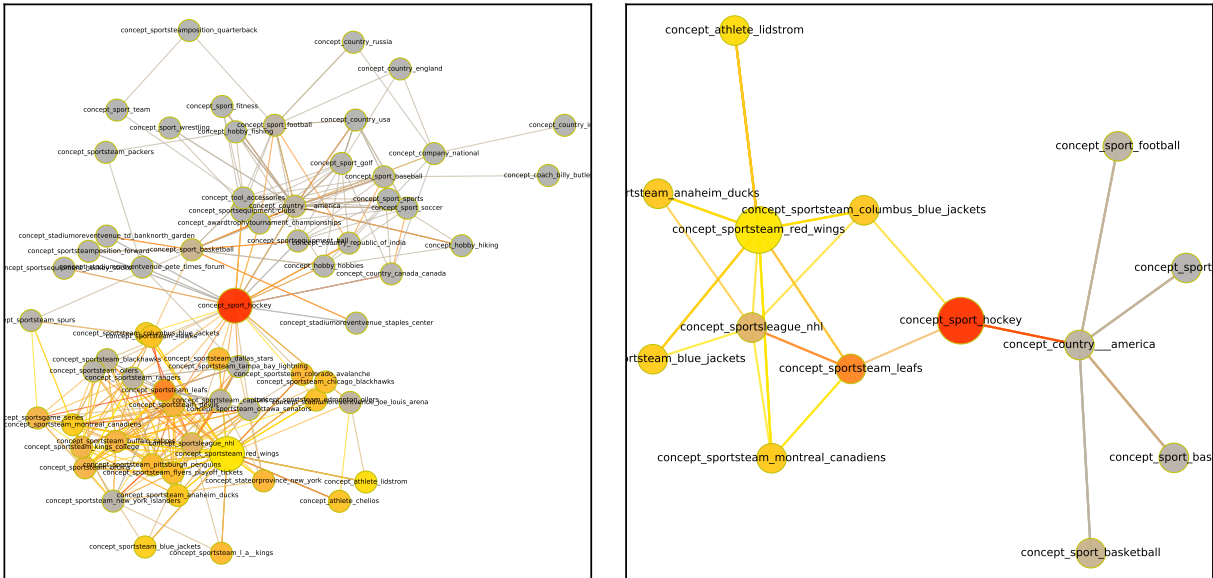


Figure 13: **TeamPlaysSport**. The head is *concept\_sportsteam\_red\_wings*, the query relation is *concept:teamplayssport*, and the tail is *concept\_sport\_hockey*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a *T*-step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the TeamPlaysSport task

Query: (concept\_sportsteam\_red\_wings , concept:teamplayssport , concept\_sport\_hockey)

Selected key edges:

- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam , concept\_sportsteam\_montreal\_canadiens
- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_montreal\_canadiens
- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam , concept\_sportsteam\_blue\_jackets
- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_blue\_jackets
- concept\_sportsteam\_red\_wings , concept:worksfor\_inv , concept\_athlete\_lidstrom
- concept\_sportsteam\_red\_wings , concept:organizationhireperson , concept\_athlete\_lidstrom
- concept\_sportsteam\_red\_wings , concept:athleteledsportsteam\_inv , concept\_athlete\_lidstrom
- concept\_sportsteam\_red\_wings , concept:athleteledsportsteam\_inv , concept\_athlete\_lidstrom
- concept\_sportsteam\_montreal\_canadiens , concept:teamplaysagainstteam , concept\_sportsteam\_red\_wings
- concept\_sportsteam\_montreal\_canadiens , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_red\_wings
- concept\_sportsteam\_montreal\_canadiens , concept:teamplaysinleague , concept\_sportsleague\_nhl
- concept\_sportsteam\_montreal\_canadiens , concept:teamplaysagainstteam , concept\_sportsteam\_leafs
- concept\_sportsteam\_montreal\_canadiens , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_leafs
- concept\_sportsteam\_blue\_jackets , concept:teamplaysagainstteam , concept\_sportsteam\_red\_wings
- concept\_sportsteam\_blue\_jackets , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_red\_wings
- concept\_sportsteam\_blue\_jackets , concept:teamplaysinleague , concept\_sportsleague\_nhl
- concept\_athlete\_lidstrom , concept:worksfor , concept\_sportsteam\_red\_wings
- concept\_athlete\_lidstrom , concept:organizationhireperson\_inv , concept\_sportsteam\_red\_wings
- concept\_athlete\_lidstrom , concept:athleteledsportsteam , concept\_sportsteam\_red\_wings
- concept\_athlete\_lidstrom , concept:athleteledsportsteam , concept\_sportsteam\_red\_wings
- concept\_sportsteam\_red\_wings , concept:teamplaysinleague , concept\_sportsleague\_nhl
- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam , concept\_sportsteam\_leafs
- concept\_sportsteam\_red\_wings , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_leafs
- concept\_sportsleague\_nhl , concept:agentcompeteswithagent , concept\_sportsleague\_nhl
- concept\_sportsleague\_nhl , concept:agentcompeteswithagent\_inv , concept\_sportsleague\_nhl
- concept\_sportsleague\_nhl , concept:teamplaysinleague\_inv , concept\_sportsteam\_leafs
- concept\_sportsteam\_leafs , concept:teamplaysinleague , concept\_sportsleague\_nhl
- concept\_sportsteam\_leafs , concept:teamplayssport , concept\_sport\_hockey

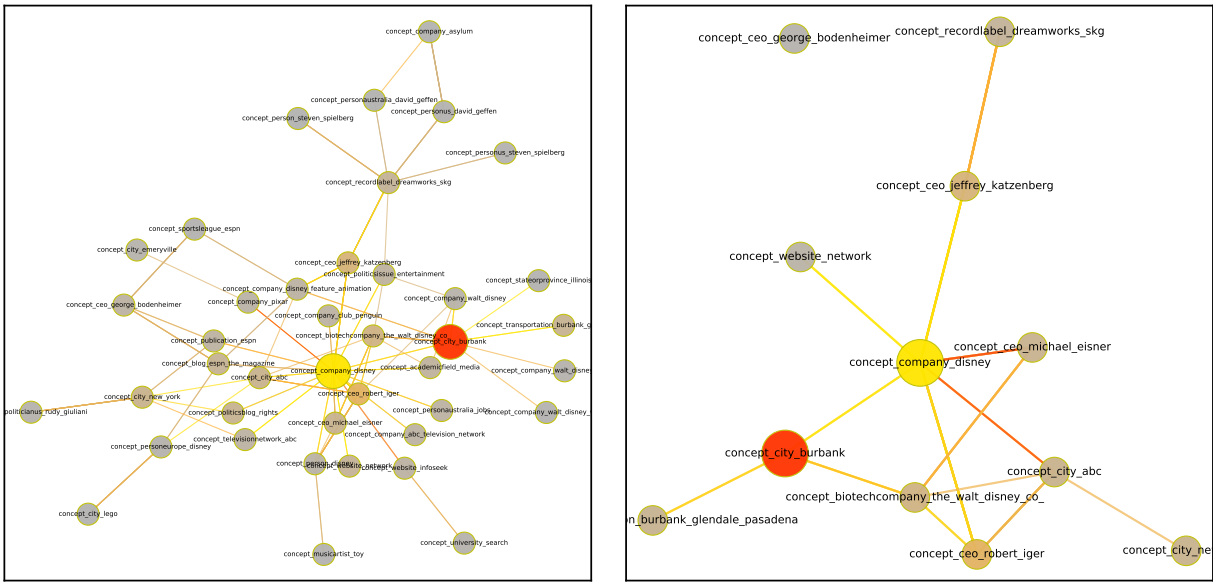


Figure 14: **OrganizationHeadQuarteredInCity**. The head is `concept_company_disney`, the query relation is `concept:organizationheadquarteredincity`, and the tail is `concept_city_burbank`. The left is a full subgraph derived with `max_attending_from_per_step=20`, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the OrganizationHeadQuarteredInCity task

Query: (`concept_company_disney`, `concept:organizationheadquarteredincity`, `concept_city_burbank`)

Selected key edges:

`concept_company_disney`, `concept:headquarteredin`, `concept_city_burbank`  
`concept_company_disney`, `concept:subpartoforganization_inv`, `concept_website_network`  
`concept_company_disney`, `concept:worksfor_inv`, `concept_ceo_robert_iger`  
`concept_company_disney`, `concept:proxyfor_inv`, `concept_ceo_robert_iger`  
`concept_company_disney`, `concept:personleadsorganization_inv`, `concept_ceo_robert_iger`  
`concept_company_disney`, `concept:ceof_inv`, `concept_ceo_robert_iger`  
`concept_company_disney`, `concept:personleadsorganization_inv`, `concept_ceo_jeffrey_katzenberg`  
`concept_company_disney`, `concept:organizationhireddperson`, `concept_ceo_jeffrey_katzenberg`  
`concept_company_disney`, `concept:organizationterminatedperson`, `concept_ceo_jeffrey_katzenberg`  
`concept_city_burbank`, `concept:headquarteredin_inv`, `concept_company_disney`  
`concept_city_burbank`, `concept:headquarteredin_inv`, `concept_biotechcompany_the_walt_disney_co_`  
`concept_website_network`, `concept:subpartoforganization`, `concept_company_disney`  
`concept_ceo_robert_iger`, `concept:worksfor`, `concept_company_disney`  
`concept_ceo_robert_iger`, `concept:proxyfor`, `concept_company_disney`  
`concept_ceo_robert_iger`, `concept:personleadsorganization`, `concept_company_disney`  
`concept_ceo_robert_iger`, `concept:ceof`, `concept_company_disney`  
`concept_ceo_robert_iger`, `concept:topmemberoforganization`, `concept_biotechcompany_the_walt_disney_co_`  
`concept_ceo_robert_iger`, `concept:organizationterminatedperson_inv`, `concept_biotechcompany_the_walt_disney_co_`  
`concept_ceo_jeffrey_katzenberg`, `concept:personleadsorganization`, `concept_company_disney`  
`concept_ceo_jeffrey_katzenberg`, `concept:organizationhireddperson_inv`, `concept_company_disney`  
`concept_ceo_jeffrey_katzenberg`, `concept:organizationterminatedperson_inv`, `concept_company_disney`  
`concept_ceo_jeffrey_katzenberg`, `concept:worksfor`, `concept_recordlabel_dreamworks_skg`  
`concept_ceo_jeffrey_katzenberg`, `concept:topmemberoforganization`, `concept_recordlabel_dreamworks_skg`  
`concept_ceo_jeffrey_katzenberg`, `concept:organizationterminatedperson_inv`, `concept_recordlabel_dreamworks_skg`  
`concept_ceo_jeffrey_katzenberg`, `concept:ceof_inv`, `concept_recordlabel_dreamworks_skg`  
`concept_biotechcompany_the_walt_disney_co_`, `concept:headquarteredin`, `concept_city_burbank`  
`concept_biotechcompany_the_walt_disney_co_`, `concept:organizationheadquarteredincity`, `concept_city_burbank`  
`concept_recordlabel_dreamworks_skg`, `concept:worksfor_inv`, `concept_ceo_jeffrey_katzenberg`  
`concept_recordlabel_dreamworks_skg`, `concept:topmemberoforganization_inv`, `concept_ceo_jeffrey_katzenberg`  
`concept_recordlabel_dreamworks_skg`, `concept:organizationterminatedperson`, `concept_ceo_jeffrey_katzenberg`  
`concept_recordlabel_dreamworks_skg`, `concept:ceof_inv`, `concept_ceo_jeffrey_katzenberg`  
`concept_city_burbank`, `concept:airportincity_inv`, `concept_transportation_burbank_glendale_pasadena`  
`concept_transportation_burbank_glendale_pasadena`, `concept:airportincity`, `concept_city_burbank`

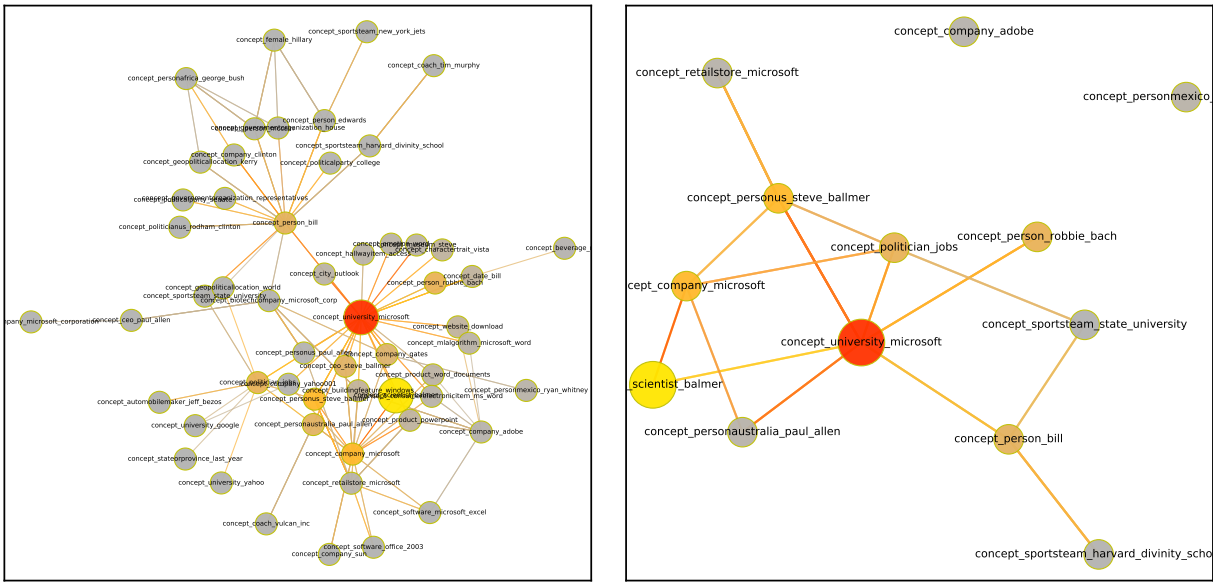


Figure 15: **WorksFor**. The head is `concept_scientist_balmer`, the query relation is `concept:worksfor`, and the tail is `concept_university_microsoft`. The left is a full subgraph derived with `max_attending_from_per_step=20`, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the WorksFor task

Query: (`concept_scientist_balmer`, `concept:worksfor`, `concept_university_microsoft`)

Selected key edges:

`concept_scientist_balmer`, `concept:topmemberoforganization`, `concept_company_microsoft`  
`concept_scientist_balmer`, `concept:organizationterminatedperson_inv`, `concept_university_microsoft`  
`concept_company_microsoft`, `concept:topmemberoforganization_inv`, `concept_personus_steve_ballmer`  
`concept_company_microsoft`, `concept:topmemberoforganization_inv`, `concept_scientist_balmer`  
`concept_university_microsoft`, `concept:agentcollaborateswithagent`, `concept_personus_steve_ballmer`  
`concept_university_microsoft`, `concept:personleadsorganization_inv`, `concept_personus_steve_ballmer`  
`concept_university_microsoft`, `concept:personleadsorganization_inv`, `concept_person_bill`  
`concept_university_microsoft`, `concept:organizationterminatedperson`, `concept_scientist_balmer`  
`concept_university_microsoft`, `concept:personleadsorganization_inv`, `concept_person_robbie_bach`  
`concept_personus_steve_ballmer`, `concept:topmemberoforganization`, `concept_company_microsoft`  
`concept_personus_steve_ballmer`, `concept:agentcollaborateswithagent_inv`, `concept_university_microsoft`  
`concept_personus_steve_ballmer`, `concept:personleadsorganization`, `concept_university_microsoft`  
`concept_personus_steve_ballmer`, `concept:worksfor`, `concept_university_microsoft`  
`concept_personus_steve_ballmer`, `concept:proxyfor`, `concept_retailstore_microsoft`  
`concept_personus_steve_ballmer`, `concept:subpartof`, `concept_retailstore_microsoft`  
`concept_personus_steve_ballmer`, `concept:agentcontrols`, `concept_retailstore_microsoft`  
`concept_person_bill`, `concept:personleadsorganization`, `concept_university_microsoft`  
`concept_person_bill`, `concept:worksfor`, `concept_university_microsoft`  
`concept_person_robbie_bach`, `concept:personleadsorganization`, `concept_university_microsoft`  
`concept_person_robbie_bach`, `concept:worksfor`, `concept_university_microsoft`  
`concept_retailstore_microsoft`, `concept:proxyfor_inv`, `concept_personus_steve_ballmer`  
`concept_retailstore_microsoft`, `concept:subpartof_inv`, `concept_personus_steve_ballmer`  
`concept_retailstore_microsoft`, `concept:agentcontrols_inv`, `concept_personus_steve_ballmer`

### For the PersonBornInLocation task

Query: (`concept_person_mark001`, `concept:personborninlocation`, `concept_county_york_city`)

Selected key edges:

`concept_person_mark001`, `concept:persongraduatedfromuniversity`, `concept_university_college`  
`concept_person_mark001`, `concept:persongraduatedschool`, `concept_university_college`  
`concept_person_mark001`, `concept:persongraduatedfromuniversity`, `concept_university_state_university`  
`concept_person_mark001`, `concept:persongraduatedschool`, `concept_university_state_university`  
`concept_person_mark001`, `concept:personbornincounty`, `concept_city_hampshire`



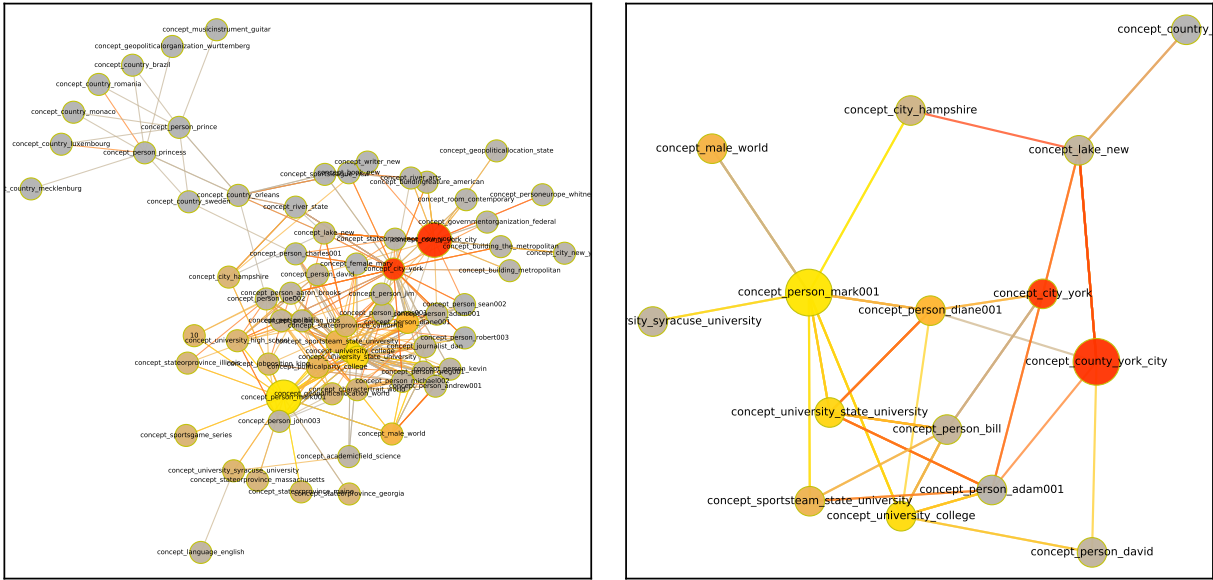


Figure 16: **PersonBornInLocation**. The head is *concept\_person\_mark001*, the query relation is *concept:personborninlocation*, and the tail is *concept\_county\_york\_city*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

```

concept_person_mark001, concept:hasspouse, concept_person_diane001
concept_person_mark001, concept:hasspouse_inv, concept_person_diane001
concept_university_college, concept:persongraduatedfromuniversity_inv, concept_person_mark001
concept_university_college, concept:persongraduatedschool_inv, concept_person_mark001
concept_university_college, concept:persongraduatedfromuniversity_inv, concept_person_bill
concept_university_college, concept:persongraduatedschool_inv, concept_person_bill
concept_university_state_university, concept:persongraduatedfromuniversity_inv, concept_person_mark001
concept_university_state_university, concept:persongraduatedschool_inv, concept_person_mark001
concept_university_state_university, concept:persongraduatedfromuniversity_inv, concept_person_bill
concept_university_state_university, concept:persongraduatedschool_inv, concept_person_bill
concept_city_hampshire, concept:personbornincity_inv, concept_person_mark001
concept_person_diane001, concept:persongraduatedfromuniversity, concept_university_state_university
concept_person_diane001, concept:persongraduatedschool, concept_university_state_university
concept_person_mark001, concept:hasspouse, concept_person_mark001
concept_person_diane001, concept:hasspouse_inv, concept_person_mark001
concept_person_diane001, concept:personborninlocation, concept_county_york_city
concept_university_state_university, concept:persongraduatedfromuniversity_inv, concept_person_diane001
concept_university_state_university, concept:persongraduatedschool_inv, concept_person_diane001
concept_person_bill, concept:personbornincity, concept_city_york
concept_person_bill, concept:personborninlocation, concept_city_york
concept_person_bill, concept:persongraduatedfromuniversity, concept_university_college
concept_person_bill, concept:persongraduatedschool, concept_university_college
concept_person_bill, concept:persongraduatedfromuniversity, concept_university_state_university
concept_person_bill, concept:persongraduatedschool, concept_university_state_university
concept_city_york, concept:personbornincity_inv, concept_person_bill
concept_city_york, concept:personbornincity_inv, concept_person_diane001
concept_university_college, concept:persongraduatedfromuniversity_inv, concept_person_diane001
concept_person_diane001, concept:personbornincity, concept_city_york
    
```

**For the PersonLeadsOrganization task**

Query: (concept\_journalist\_bill\_plante, concept:personleadsorganization, concept\_company\_cnn\_pbs)

Selected key edges:  
concept\_journalist\_bill\_plante, concept:worksfor, concept\_televisionnetwork\_cbs  
concept\_journalist\_bill\_plante, concept:agentcollaborateswithagent\_inv, concept\_televisionnetwork\_cbs  
concept\_televisionnetwork\_cbs, concept:worksfor\_inv, concept\_journalist\_walter\_cronkite

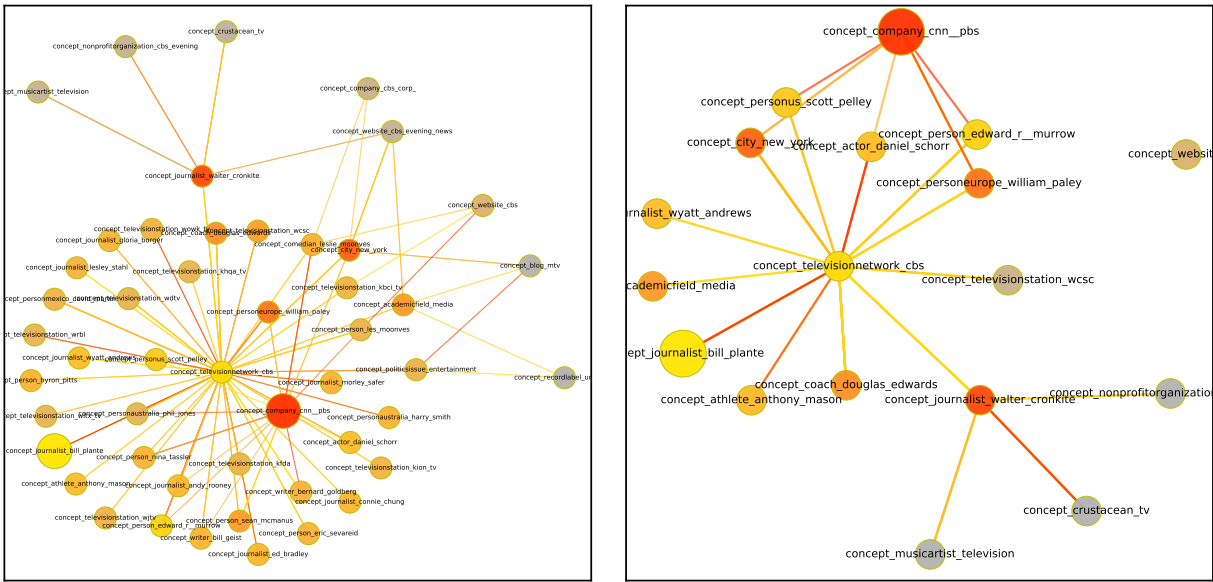


Figure 17: **PersonLeadsOrganization**. The head is `concept_journalist_bill_plante`, the query relation is `concept:organizationheadquarteredincity`, and the tail is `concept_company_cnn_pbs`. The left is a full subgraph derived with `max_attending_from_per_step=20`, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

`concept_televisionnetwork_cbs`, `concept:agentcollaborateswithagent`, `concept_journalist_walter_cronkite`  
`concept_televisionnetwork_cbs`, `concept:worksfor_inv`, `concept_personus_scott_pelley`  
`concept_televisionnetwork_cbs`, `concept:worksfor_inv`, `concept_actor_daniel_schorr`  
`concept_televisionnetwork_cbs`, `concept:worksfor_inv`, `concept_person_edward_r__murrow`  
`concept_televisionnetwork_cbs`, `concept:agentcollaborateswithagent`, `concept_person_edward_r__murrow`  
`concept_televisionnetwork_cbs`, `concept:worksfor_inv`, `concept_journalist_bill_plante`  
`concept_televisionnetwork_cbs`, `concept:agentcollaborateswithagent`, `concept_journalist_bill_plante`  
`concept_journalist_walter_cronkite`, `concept:worksfor`, `concept_televisionnetwork_cbs`  
`concept_journalist_walter_cronkite`, `concept:agentcollaborateswithagent_inv`, `concept_televisionnetwork_cbs`  
`concept_journalist_walter_cronkite`, `concept:worksfor`, `concept_nonprofitorganization_cbs_evening`  
`concept_personus_scott_pelley`, `concept:worksfor`, `concept_televisionnetwork_cbs`  
`concept_personus_scott_pelley`, `concept:personleadsorganization`, `concept_televisionnetwork_cbs`  
`concept_personus_scott_pelley`, `concept:personleadsorganization`, `concept_company_cnn_pbs`  
`concept_actor_daniel_schorr`, `concept:worksfor`, `concept_televisionnetwork_cbs`  
`concept_actor_daniel_schorr`, `concept:personleadsorganization`, `concept_televisionnetwork_cbs`  
`concept_actor_daniel_schorr`, `concept:personleadsorganization`, `concept_company_cnn_pbs`  
`concept_person_edward_r__murrow`, `concept:worksfor`, `concept_televisionnetwork_cbs`  
`concept_person_edward_r__murrow`, `concept:agentcollaborateswithagent_inv`, `concept_televisionnetwork_cbs`  
`concept_person_edward_r__murrow`, `concept:personleadsorganization`, `concept_televisionnetwork_cbs`  
`concept_person_edward_r__murrow`, `concept:personleadsorganization`, `concept_company_cnn_pbs`  
`concept_televisionnetwork_cbs`, `concept:organizationheadquarteredincity`, `concept_city_new_york`  
`concept_televisionnetwork_cbs`, `concept:headquarteredin`, `concept_city_new_york`  
`concept_televisionnetwork_cbs`, `concept:agentcollaborateswithagent`, `concept_person_europe_william_paley`  
`concept_televisionnetwork_cbs`, `concept:topmemberoforganization_inv`, `concept_person_europe_william_paley`  
`concept_company_cnn_pbs`, `concept:headquarteredin`, `concept_city_new_york`  
`concept_company_cnn_pbs`, `concept:personbelongstoorganization_inv`, `concept_person_europe_william_paley`  
`concept_nonprofitorganization_cbs_evening`, `concept:worksfor_inv`, `concept_journalist_walter_cronkite`  
`concept_city_new_york`, `concept:organizationheadquarteredincity_inv`, `concept_televisionnetwork_cbs`  
`concept_city_new_york`, `concept:headquarteredin_inv`, `concept_televisionnetwork_cbs`  
`concept_city_new_york`, `concept:headquarteredin_inv`, `concept_company_cnn_pbs`  
`concept_person_europe_william_paley`, `concept:agentcollaborateswithagent_inv`, `concept_televisionnetwork_cbs`  
`concept_person_europe_william_paley`, `concept:topmemberoforganization`, `concept_televisionnetwork_cbs`  
`concept_person_europe_william_paley`, `concept:personbelongstoorganization`, `concept_company_cnn_pbs`  
`concept_person_europe_william_paley`, `concept:personleadsorganization`, `concept_company_cnn_pbs`

**For the OrganizationHiredPerson task**

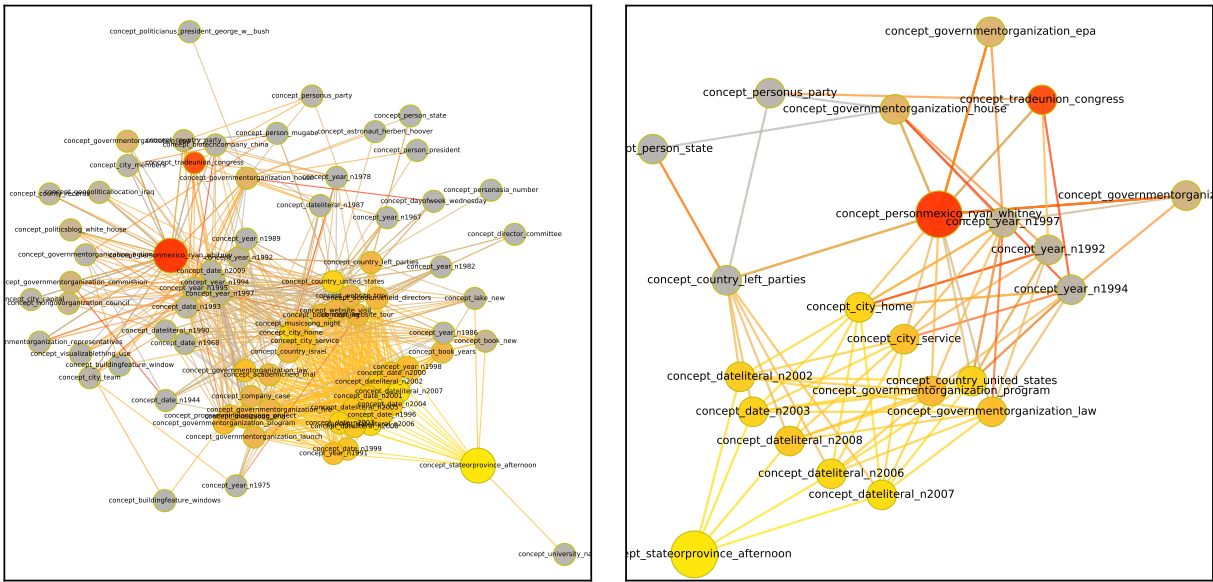


Figure 18: **OrganizationHiredPerson**. The head is *concept\_stateorprovince\_afternoon*, the query relation is *concept:organizationhiredperson*, and the tail is *concept\_personmexico-ryan\_whitney*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

Query: (concept.stateorprovince\_afternoon, concept:organizationhiredperson, concept.personmexico-ryan\_whitney)

Selected key edges:

concept.stateorprovince\_afternoon, concept:atdate, concept.dateliteral\_n2007  
concept.stateorprovince\_afternoon, concept:atdate, concept.date\_n2003  
concept.stateorprovince\_afternoon, concept:atdate, concept.dateliteral\_n2006  
concept.dateliteral\_n2007, concept:atdate\_inv, concept.country\_united\_states  
concept.dateliteral\_n2007, concept:atdate\_inv, concept.city\_home  
concept.dateliteral\_n2007, concept:atdate\_inv, concept.city\_service  
concept.dateliteral\_n2007, concept:atdate\_inv, concept.country\_left\_parties  
concept.date\_n2003, concept:atdate\_inv, concept.country\_united\_states  
concept.date\_n2003, concept:atdate\_inv, concept.city\_home  
concept.date\_n2003, concept:atdate\_inv, concept.city\_service  
concept.date\_n2003, concept:atdate\_inv, concept.country\_left\_parties  
concept.dateliteral\_n2006, concept:atdate\_inv, concept.country\_united\_states  
concept.dateliteral\_n2006, concept:atdate\_inv, concept.city\_home  
concept.dateliteral\_n2006, concept:atdate\_inv, concept.city\_service  
concept.dateliteral\_n2006, concept:atdate\_inv, concept.country\_left\_parties  
concept.country\_united\_states, concept:atdate, concept.year\_n1992  
concept.country\_united\_states, concept:atdate, concept.year\_n1997  
concept.country\_united\_states, concept:organizationhiredperson, concept.personmexico-ryan\_whitney  
concept.city\_home, concept:atdate, concept.year\_n1992  
concept.city\_home, concept:atdate, concept.year\_n1997  
concept.city\_home, concept:organizationhiredperson, concept.personmexico-ryan\_whitney  
concept.city\_service, concept:atdate, concept.year\_n1992  
concept.city\_service, concept:atdate, concept.year\_n1997  
concept.city\_service, concept:organizationhiredperson, concept.personmexico-ryan\_whitney  
concept.country\_left\_parties, concept:worksfor\_inv, concept.personmexico-ryan\_whitney  
concept.country\_left\_parties, concept:organizationhiredperson, concept.personmexico-ryan\_whitney  
concept.year\_n1992, concept:atdate\_inv, concept.governmentorganization\_house  
concept.year\_n1992, concept:atdate\_inv, concept.country\_united\_states  
concept.year\_n1992, concept:atdate\_inv, concept.city\_home  
concept.year\_n1992, concept:atdate\_inv, concept.tradeunion\_congress  
concept.year\_n1997, concept:atdate\_inv, concept.governmentorganization\_house  
concept.year\_n1997, concept:atdate\_inv, concept.country\_united\_states  
concept.year\_n1997, concept:atdate\_inv, concept.city\_home  
concept.personmexico-ryan\_whitney, concept:worksfor, concept.governmentorganization\_house

concept\_personmexico\_ryan\_whitney, concept:worksfor, concept:tradeunion\_congress  
 concept\_personmexico\_ryan\_whitney, concept:worksfor, concept\_country\_left\_parties  
 concept\_governmentorganization\_house, concept:personbelongstoorganization\_inv, concept\_personus\_party  
 concept\_governmentorganization\_house, concept:worksfor\_inv, concept\_personmexico\_ryan\_whitney  
 concept\_governmentorganization\_house, concept:organizationhireperson, concept\_personmexico\_ryan\_whitney  
 concept\_tradeunion\_congress, concept:organizationhireperson, concept\_personus\_party  
 concept\_tradeunion\_congress, concept:worksfor\_inv, concept\_personmexico\_ryan\_whitney  
 concept\_tradeunion\_congress, concept:organizationhireperson, concept\_personmexico\_ryan\_whitney  
 concept\_country\_left\_parties, concept:organizationhireperson, concept\_personus\_party

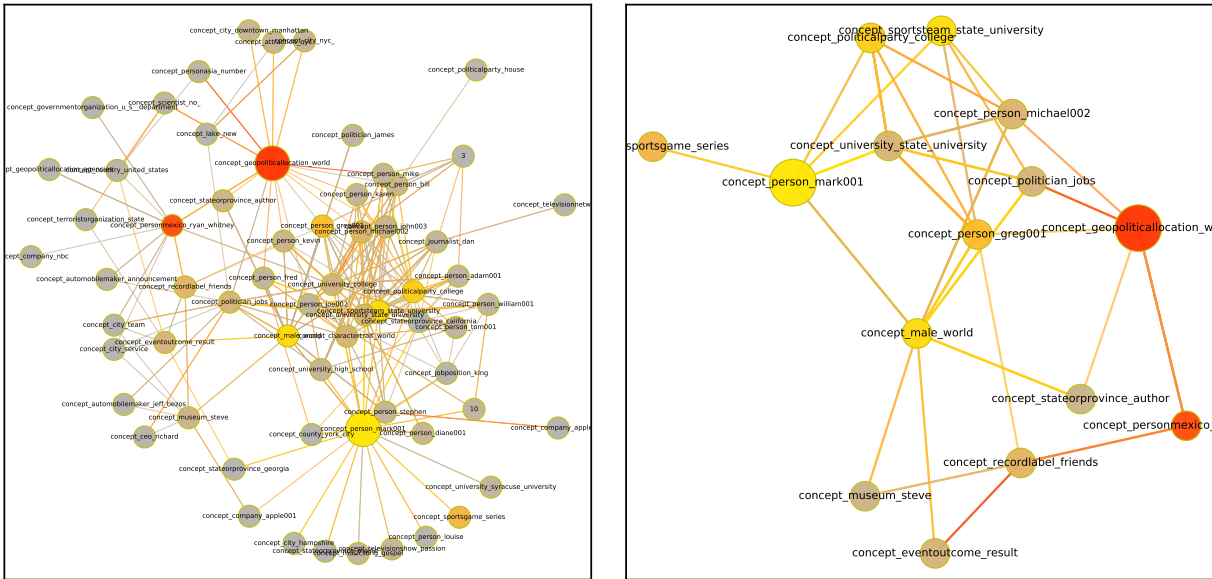


Figure 19: **AgentBelongsToOrganization**. The head is *concept\_person\_mark001*, the query relation is *concept:agentbelongstoorganization*, and the tail is *concept\_geopoliticallocation\_world*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

**For the AgentBelongsToOrganization task**

Query: (concept\_person\_mark001, concept:agentbelongstoorganization, concept\_geopoliticallocation\_world)

Selected key edges:

- concept\_person\_mark001, concept:personbelongstoorganization, concept\_sportsteam\_state\_university
- concept\_person\_mark001, concept:agentcollaborateswithagent, concept\_male\_world
- concept\_person\_mark001, concept:agentcollaborateswithagent\_inv, concept\_male\_world
- concept\_person\_mark001, concept:personbelongstoorganization, concept\_politicalparty\_college
- concept\_sportsteam\_state\_university, concept:personbelongstoorganization\_inv, concept\_politician\_jobs
- concept\_sportsteam\_state\_university, concept:personbelongstoorganization\_inv, concept\_person\_mark001
- concept\_sportsteam\_state\_university, concept:personbelongstoorganization\_inv, concept\_person\_greg001
- concept\_sportsteam\_state\_university, concept:personbelongstoorganization\_inv, concept\_person\_michael002
- concept\_male\_world, concept:agentcollaborateswithagent, concept\_politician\_jobs
- concept\_male\_world, concept:agentcollaborateswithagent\_inv, concept\_politician\_jobs
- concept\_male\_world, concept:agentcollaborateswithagent, concept\_person\_mark001
- concept\_male\_world, concept:agentcollaborateswithagent\_inv, concept\_person\_mark001
- concept\_male\_world, concept:agentcollaborateswithagent, concept\_person\_greg001
- concept\_male\_world, concept:agentcollaborateswithagent\_inv, concept\_person\_greg001
- concept\_male\_world, concept:agentcontrols, concept\_person\_greg001
- concept\_male\_world, concept:agentcollaborateswithagent, concept\_person\_michael002
- concept\_male\_world, concept:agentcollaborateswithagent\_inv, concept\_person\_michael002
- concept\_politicalparty\_college, concept:personbelongstoorganization\_inv, concept\_person\_mark001
- concept\_politicalparty\_college, concept:personbelongstoorganization\_inv, concept\_person\_greg001
- concept\_politicalparty\_college, concept:personbelongstoorganization\_inv, concept\_person\_michael002
- concept\_politician\_jobs, concept:personbelongstoorganization, concept\_sportsteam\_state\_university
- concept\_politician\_jobs, concept:agentcollaborateswithagent, concept\_male\_world

concept\_politician\_jobs , concept\_agentcollaborateswithagent\_inv , concept\_male\_world  
concept\_politician\_jobs , concept\_worksfor , concept\_geopoliticallocation\_world  
concept\_person\_greg001 , concept\_personbelongstoorganization , concept\_sportsteam\_state\_university  
concept\_person\_greg001 , concept\_agentcollaborateswithagent , concept\_male\_world  
concept\_person\_greg001 , concept\_agentcollaborateswithagent\_inv , concept\_male\_world  
concept\_person\_greg001 , concept\_agentcontrols\_inv , concept\_male\_world  
concept\_person\_greg001 , concept\_agentbelongstoorganization , concept\_geopoliticallocation\_world  
concept\_person\_greg001 , concept\_personbelongstoorganization , concept\_politicalparty\_college  
concept\_person\_greg001 , concept\_agentbelongstoorganization , concept\_recordlabel\_friends  
concept\_person\_michael002 , concept\_personbelongstoorganization , concept\_sportsteam\_state\_university  
concept\_person\_michael002 , concept\_agentcollaborateswithagent , concept\_male\_world  
concept\_person\_michael002 , concept\_agentcollaborateswithagent\_inv , concept\_male\_world  
concept\_person\_michael002 , concept\_agentbelongstoorganization , concept\_geopoliticallocation\_world  
concept\_person\_michael002 , concept\_personbelongstoorganization , concept\_politicalparty\_college  
concept\_geopoliticallocation\_world , concept\_worksfor\_inv , concept\_personmexico\_ryan\_whitney  
concept\_geopoliticallocation\_world , concept\_organizationhireperson , concept\_personmexico\_ryan\_whitney  
concept\_recordlabel\_friends , concept\_organizationhireperson , concept\_personmexico\_ryan\_whitney  
concept\_personmexico\_ryan\_whitney , concept\_worksfor , concept\_geopoliticallocation\_world  
concept\_personmexico\_ryan\_whitney , concept\_organizationhireperson\_inv , concept\_geopoliticallocation\_world  
concept\_personmexico\_ryan\_whitney , concept\_organizationhireperson\_inv , concept\_recordlabel\_friends

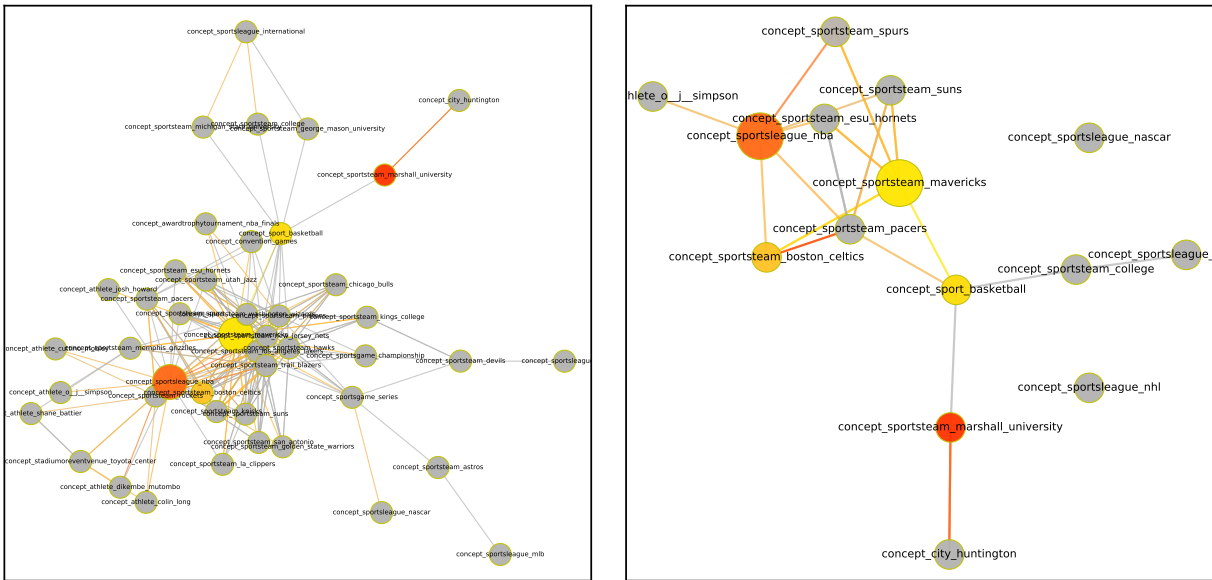


Figure 20: **TeamPlaysInLeague**. The head is *concept\_sportsteam\_mavericks*, the query relation is *concept:teamplaysinleague*, and the tail is *concept\_sportsleague\_nba*. The left is a full subgraph derived with *max\_attending\_from\_per\_step=20*, and the right is a further pruned subgraph from the left based on attention. The big yellow node represents the head, and the big red node represents the tail. Color on the rest indicates attention scores over a  $T$ -step reasoning process, where grey means less attention, yellow means more attention gained during early steps, and red means gaining more attention when getting closer to the final step.

### For the TeamPlaysInLeague task

Query: (concept\_sportsteam\_mavericks , concept:teamplaysinleague , concept\_sportsleague\_nba)

Selected key edges:

concept\_sportsteam\_mavericks , concept:teamplayssport , concept\_sport\_basketball  
concept\_sportsteam\_mavericks , concept:teamplaysagainstteam , concept\_sportsteam\_boston\_celtics  
concept\_sportsteam\_mavericks , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_boston\_celtics  
concept\_sportsteam\_mavericks , concept:teamplaysagainstteam , concept\_sportsteam\_spurs  
concept\_sportsteam\_mavericks , concept:teamplaysagainstteam\_inv , concept\_sportsteam\_spurs  
concept\_sport\_basketball , concept:teamplayssport\_inv , concept\_sportsteam\_college  
concept\_sport\_basketball , concept:teamplayssport\_inv , concept\_sportsteam\_marshall\_university  
concept\_sportsteam\_boston\_celtics , concept:teamplaysinleague , concept\_sportsleague\_nba  
concept\_sportsteam\_spurs , concept:teamplaysinleague , concept\_sportsleague\_nba  
concept\_sportsleague\_nba , concept:agentcompeteswithagent , concept\_sportsleague\_nba

concept:sportsleague\_nba , concept:agentcompeteswithagent\_inv , concept:sportsleague\_nba  
concept:sportsteam\_college , concept:teamplaysinleague , concept:sportsleague\_international