

# Efficient Link Prediction via GNN Layers Induced by Negative Sampling

Yuxin Wang, Xiannian Hu, Quan Gan, Xuanjing Huang, Xipeng Qiu, David Wipf *Fellow, IEEE*

**Abstract**—<sup>1</sup>Graph neural networks (GNNs) for link prediction can loosely be divided into two broad categories. First, *node-wise* architectures pre-compute individual embeddings for each node that are later combined by a simple decoder to make predictions. While extremely efficient at inference time, model expressiveness is limited such that isomorphic nodes contributing to candidate edges may not be distinguishable, compromising accuracy. In contrast, *edge-wise* methods rely on the formation of edge-specific subgraph embeddings to enrich the representation of pair-wise relationships, disambiguating isomorphic nodes to improve accuracy, but with increased model complexity. To better navigate this trade-off, we propose a novel GNN architecture whereby the *forward pass* explicitly depends on *both* positive (as is typical) and negative (unique to our approach) edges to inform more flexible, yet still cheap node-wise embeddings. This is achieved by recasting the embeddings themselves as minimizers of a forward-pass-specific energy function that favors separation of positive and negative samples. Notably, this energy is distinct from the actual training loss shared by most existing link prediction models, where contrastive pairs only influence the *backward pass*. As demonstrated by extensive empirical evaluations, the resulting architecture retains the inference speed of node-wise models, while producing competitive accuracy with edge-wise alternatives. We released our code at <https://github.com/yxzwang/SubmissionverOfYinYanGNN>.

**Index Terms**—Artificial Intelligence, Machine Learning, Graph Neural Networks, Link Prediction

## I. INTRODUCTION

**L**INK Prediction is a fundamental graph learning challenge that involves determining whether or not there should exist an edge connecting two nodes. Given the prevalence of graph-structured data, this task has widespread practical significance spanning domains such as social networks, knowledge graphs, and e-commerce, and recommendation systems [1], [2], [3]. As one representative example of the latter, the goal could be to predict whether a user node should be linked with an item node in a product graph, where edges are indicative of some form of user/item engagement, e.g., clicks, purchases, etc.

Beyond heuristic techniques such as Common Neighbors (CN) [4], Adamic-Adar (AA) [5], and Resource Allocation (RA) [6], graph neural networks (GNNs) have recently shown tremendous promise in addressing link prediction with trainable deep architectures [7], [8], [9], [10], [11], [12]. Broadly speaking these GNN models fall into two categories, based on whether they rely on *node-wise* or *edge-wise* embeddings.

The former involves using a GNN to pre-compute individual embeddings for each node that are later combined by a simple decoder to predict the presence of edges. This strategy is preferable when inference speed is paramount (as is often the case in real-world applications requiring low-latency predictions), since once node-wise embeddings are available, combining them to make predictions is cheap. Moreover, accelerated decoding techniques such as Maximum Inner Product Search (MIPS) [13], [14], [15] or Flashlight [16] exist to further economize inference. The downside though of node-wise embedding methods is that they may fail to disambiguate isomorphic nodes that combine to form a candidate edge [9].

To this end, edge-wise embeddings with greater expressive power have been proposed for more robust link prediction [10], [17], [12], [18]. These models base their predictions on edge-specific subgraphs capable of breaking isomorphic node relationships via structural information (e.g., overlapping neighbors, shortest paths, positional encodings, or subgraph sketches) that might otherwise undermine the performance of node-wise embeddings. This flexibility comes with a substantial cost though, as inference complexity can be orders of magnitude larger given that a unique subgraph must be extracted and processed by the GNN for every test edge. Although this expense can be alleviated to some cases by pre-processing [12], for inference over very large sets of candidate links, even the pre-processing time can be overwhelming relative to that required by node-wise predictors.

In this work, we address the trade-off between expressiveness and inference efficiency via the following strategy. To maintain minimal inference speeds, we restrict ourselves to a node-wise embedding approach and then try to increase the expressiveness on multiple fronts. Most importantly, we allow each node-level embedding computed during the forward pass to depend on not only its ego network (i.e., subgraph containing a target node), but also on the embeddings of negatively sampled nodes, meaning nodes that were not originally sharing an edge with the target node. This can be viewed as forming a complementary *negative* ego network for each node. Moreover, rather than heuristically incorporating the resulting positive *and* negative ego networks within a traditional GNN-based embedding model, we instead combine them so as to infuse their integration with an inductive bias specifically tailored for link prediction. Specifically, we introduce a parameterized graph-regularized energy, in the spirit of triplet ranking loss functions used for capturing both relative similarities and differences between pair-wise items. By design, the parameter-dependent minimizer of this function can then serve the role of end-to-end trainable node-wise embeddings, *explicitly*

The authors are with the School of Computer Science, Fudan University (Yuxin Wang, Xiannian Hu, Xuanjing Huang, Xipeng Qiu) and Amazon (Quan Gan, David Wipf).

E-mail: {wangyuxin21,21210240194}@m.fudan.edu.cn, quagan@amazon.com, {xjhuang,xpqi}@fudan.edu.cn, davidwipf@gmail.com <sup>1</sup> This article has been accepted for publication in IEEE Transactions on Knowledge and Data Engineering. Citation information: DOI 10.1109/TKDE.2024.3481015

TABLE I

COMPARISONS WITH EXISTING NODE-WISE AND EDGE-WISE MODELS. OVERALL, YINYANGNN CAN ACHIEVE ACCURACY COMPARABLE TO EDGE-WISE MODELS WHILE MAINTAINING THE EFFICIENCY OF NODE-WISE MODELS. YINYANGNN ALSO POSSESSES MULTIPLE DESIRABLE PROPERTIES LINKED TO ITS UNIQUE ARCHITECTURAL DESIGN.

|  | Node-wise Models | Edge-wise Models | YinYanGNN |
|--|------------------|------------------|-----------|
| Can distinguish some isomorphic node pairs | ×                | ✓                | ✓         |
| Competitive accuracy                       | ×                | ✓                | ✓         |
| Fast inference run-time                    | ✓                | ×                | ✓         |
| Sublinear acceleration                     | ✓                | ×                | ✓         |
| Energy descent convergence guarantees      | ×                | ×                | ✓         |
| Negative samples in forward pass           | ×                | ×                | ✓         |

dependent on the node features of both positive and negative samples even during the forward pass (not just the backward training pass as is typical). For these reasons, we refer to our model as a *Yin* (negative) *Yang* (positive) GNN, or YinYanGNN for short.

In this way, we increase the flexibility of node-wise embedding approaches, without significantly increasing the computational complexity, as no edge-wise embeddings or edge-specific subgraph extraction is necessary. Additionally, by unifying the positive and negative samples within a single energy function minimization process, the implicit receptive field of the embeddings can be arbitrarily large without oversmoothing, a property we inherit from prior related work on optimization-based GNN models applied to much simpler node classification tasks [19]. These observations lead to a statement of our primary contributions:

- 1) We design node-wise embeddings for link prediction that are explicitly imbued with an inductive bias informed by the node features of *both* positive and negative samples during the model *forward pass*, not merely the *backward pass* as is customary with contrastive learning-based methods. This is accomplished by recasting the embeddings themselves as minimizers of an energy function that explicitly balances the impact of positive (*Yang*) and negative (*Yin*) samples, leading to a model we refer to as YinYanGNN.
- 2) We analyze the convergence properties and computational complexity of the optimization process which produces YinYanGNN embeddings, as well as their expressiveness relative to traditional node-wise models. These results suggest that our approach can potentially serve as a reliable compromise between node- and edge-wise alternatives.
- 3) Experiments on real-world link prediction benchmarks reveal the YinYanGNN can outperform SOTA node-wise models in terms of accuracy while matching their efficiency. And analogously, YinYanGNN can exceed the efficiency of edge-wise approaches while maintaining similar (and in some cases better) prediction accuracy.

We summarize the differentiating factors of YinYanGNN compared to existing node-wise and edge-wise models in Table I, while the basic architecture is displayed in Figure 1, with more detailed descriptions to follow in subsequent sections.

## II. RELATED WORK

**GNNs for Link Prediction.** As mentioned in Section I, GNN models for link prediction can be roughly divided into two categories, those based on node-wise embeddings [7], [8], [20] and those based on edge-wise embeddings [9], [12], [17], [11], [21], [22], [18]. The former is generally far more efficient at inference time given that the embeddings need only be computed once for each node and then repeatedly combined to make predictions for each candidate edge. However, the latter is more expressive by facilitating edge-specific structural features at the cost of much slower inference.

**GNN Layers formed from unfolded optimization steps.** A plethora of recent research has showcased the potential of constructing resilient GNN architectures for node classification using graph propagation layers that emulate the iterative descent steps of a graph-regularized energy function [23], [24], [25], [26], [27], [19], [28], [29]. These approaches allow the node embeddings at each layer to be regarded as progressively refined approximations of an interpretable energy minimizer. A key advantage is that embeddings obtained in this way can be purposefully designed to address challenges such as GNN oversmoothing or the introduction of robustness against spurious edges. Moreover, these adaptable embeddings can be seamlessly integrated into a bilevel optimization framework [30] for supervised training. Even so, prior work in this domain thus far has been primarily limited to much simpler node classification tasks, where nuanced relationships between pairs of nodes need not be explicitly accounted for. In contrast, we are particularly interested in the latter, and the potential to design new energy functions that introduce inductive biases suitable for link prediction.

**GNN Expressiveness.** Most work on GNN expressiveness focuses on the representational power of entire graphs[31], [32], [33], [34], [35], [36], [37]. More narrowly in terms of link prediction, structural representations of node sets have also been studied[38], as well as labeling tricks [10] first proposed by [9] to improve expressiveness. Overall, there now exists a variety of edge-wise methods[12], [18], [22] specifically designed to focus on more expressive edge-specific structural features as we discuss elsewhere herein.

**Contrastive Learning.** Although negative pairs and contrastive learning are commonly used by link prediction frameworks and self-supervised learning more generally [39], [40], [41], their role in prior work is largely restricted to defining the training loss, and therefore only influencing the backward pass. This is unlike YinYanGNN, whereby the negative samples also explicitly impact the forward pass and core model/encoder architecture itself.

## III. PRELIMINARIES

In this section we briefly introduce notation before providing concrete details of the link prediction problem that will be useful later.

### A. Notation

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, X)$  be a graph with node set  $\mathcal{V}$ , corresponding  $d_x$ -dimensional node features  $X \in \mathbb{R}^{n \times d_x}$ , and edge set  $\mathcal{E}$ ,

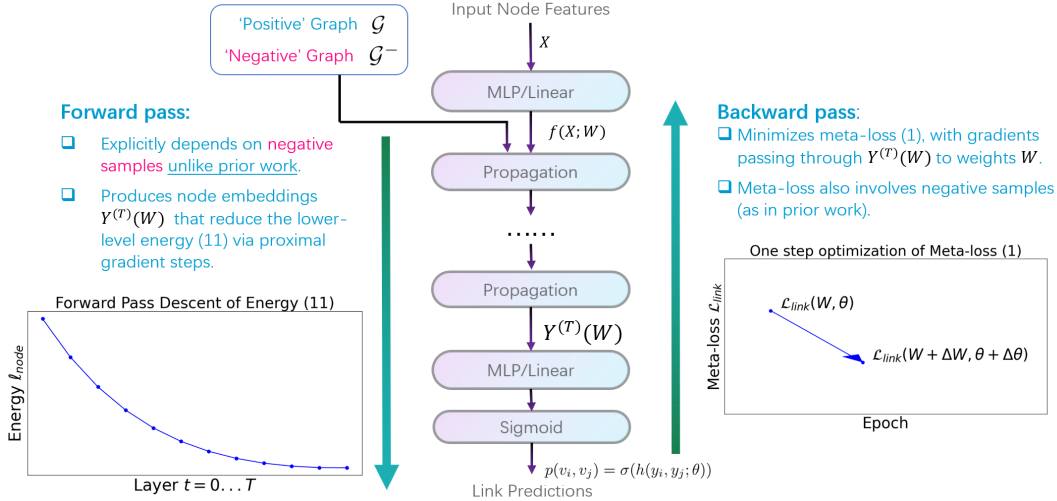


Fig. 1. *YinYanGNN model illustration.* On the left side we show the YinYanGNN forward pass explicitly depending on negative samples, with layers computing embeddings that descend a lower-level energy  $\ell_{node}$  defined by (11). On the right side we show the more traditional backward pass for optimizing meta-loss (1) and one-step optimization of it over parameters  $W$  (which define the lower-level energy) and  $\theta$  (specific to the meta-loss). Supporting details and derivations will be presented in Section IV.

where  $|\mathcal{V}| = n$ . We use  $A$  to denote the adjacency matrix and  $D$  for the degree matrix. The associated Laplacian matrix is defined by  $L \triangleq D - A$ . Furthermore,  $Y \in \mathbb{R}^{n \times d}$  refers to node embeddings of size  $d$  we seek to learn via a node-wise link prediction procedure. Specifically the node embedding for node  $i$  is  $y_i$ , which is equivalent to the  $i$ -th row of  $Y$ .

### B. Link Prediction

We begin by introducing a commonly-used loss for link prediction, which is defined over the training set  $\mathcal{E}_{train} \subset \mathcal{E}$ . For both node-wise and edge-wise methods, the shared goal is to obtain an edge probability score  $p(v_i, v_j) = \sigma(e_{ij})$  for all edges  $(v_i, v_j) \in \mathcal{E}_{train}$  (as well as negatively-sampled counterparts to be determined shortly), where  $\sigma$  is a sigmoid function and  $e_{ij}$  is a discriminative representation for edge  $(v_i, v_j)$ . Proceeding further, for every true positive edge  $(v_i, v_j)$  in the training set,  $N \geq 1$  negative edges  $(v_i^a, v_j^a)_{a=1, \dots, N}$  are randomly sampled from the graph for supervision purposes. We are then positioned to express the overall link prediction loss as  $\mathcal{L}_{link} \triangleq =$

$$\sum_{(i,j) \in \mathcal{E}_{train}} \left[ -\log(p(v_i, v_j)) - \sum_{a=1}^N \frac{1}{N} \log(1 - p(v_i^a, v_j^a)) \right], \quad (1)$$

where each edge probability is computed with the corresponding edge representation. The lingering difference between node- and edge-wise methods then lies in how each edge representation  $e_{ij}$  is actually computed.

For node-wise methods,  $e_{ij} = h(y_i, y_j)$ , where  $y_i$  and  $y_j$  are node-wise embeddings and  $h$  is a decoder function ranging in complexity from a parameter-free inner-product to a multi-layer MLP. While decoder structure varies [42], [43], [44], [45], [16], of particular note for its practical effectiveness is the HadamardMLP approach, which amounts to simply computing

the hadamard product between  $y_i$  and  $y_j$  and then passing the result through an MLP with trainable parameters  $\theta$ . Fast, sublinear inference times are possible with HadamardMLP using an algorithm from [16]. In contrast, the constituent node embeddings themselves are typically computed with some form of trainable GNN encoder model  $g$  of the form  $y_i = g(x_i, \mathcal{G}_i)$  and  $y_j = g(x_j, \mathcal{G}_j)$ , where  $\mathcal{G}_i$  and  $\mathcal{G}_j$  are the subgraphs containing nodes  $v_i$  and  $v_j$ , respectively.

Turning to edge-wise methods, the edge representation  $e_{ij}$  relies on the subgraph  $\mathcal{G}_{ij}$  defined by *both*  $v_i$  and  $v_j$ . In this case  $e_{ij} = h_e(v_i, v_j, \mathcal{G}_{ij})$ , where  $h_e$  is an edge encoder GNN whose predictions can generally *not* be decomposed into a function of individual node embeddings as before. Note also that while the embeddings from node-wise subgraphs for *all* nodes in the graph can be produced by a *single* GNN forward pass, a unique/separate edge-wise subgraph and corresponding forward pass are needed to make predictions for each candidate edge. This explains why edge-wise models endure far slower inference speeds in practice.

## IV. INCORPORATING NEGATIVE SAMPLING INTO NODE-WISE MODEL DESIGN

Previously we described how computationally-efficient node-wise embedding methods for link prediction rely on edge representations that decompose as  $e_{ij} = h[g(y_i, \mathcal{G}_i), g(y_j, \mathcal{G}_j)]$  for node-pair  $(v_i, v_j)$ , a decomposition that is decidedly less expressive than the more general form  $e_{ij} = h_e(v_i, v_j, \mathcal{G}_{ij})$  adopted by edge-wise embedding methods. Although we can never match the flexibility of the edge-wise models with a node-wise approach, we can nonetheless increase the expressiveness of node-wise models while still retaining their attractive computational footprint.

At a high-level, our strategy for accomplishing this goal is to learn node-wise embeddings of the revised form  $y_i = g(v_i, \mathcal{G}_i, \mathcal{G}_i^-)$ , where  $\mathcal{G}_i^-$  is a subgraph of  $\mathcal{G}^-$  centered at node

$v_i$ ,  $\mathcal{G}^- = (\mathcal{V}, \mathcal{E}^-, X)$ , and  $\mathcal{E}^-$  is a set of negatively-sampled edges between nodes in the original graph  $\mathcal{G}$ . In this way each node-wise embedding has access to node features from both positive and negative neighboring nodes.

To operationalize this conceptual design, rather than heuristically embedding negative samples within an existing GNN architecture (see Section VII-D for experiments using this simple strategy), we instead chose node-wise embeddings that minimize an energy function regularized by both positive and negative edges, i.e.,  $\mathcal{E}$  and  $\mathcal{E}^-$ . More formally, we seek a node embedding matrix  $Y = \arg \min_Y \ell_{node}(\mathcal{G}, \mathcal{G}^-)$  in such a way that the optimal solution decomposes as  $y_i = g(v_i, \mathcal{G}_i, \mathcal{G}_i^-)$  for some differentiable function  $g$  across all nodes  $v_i$ . This allows us to anchor the influence of positive and negative edges within a unified energy surface, with trainable minimizers that can be embedded within the link prediction loss from (1). In the remainder of this section we motivate our selection for  $\ell_{node}$ , as well as the optimization steps which form the structure of the corresponding function  $g$ .

### A. An Initial Energy Function

Prior work on optimization-based node embeddings [46], [26], [27], [19], [28], [29] largely draw on energy functions related to [47], which facilitates the balancing of local consistency relative to labels or a base predictor, with global constraints from graph structure. However, these desiderata alone are inadequate for the link prediction task, where we would also like to drive individual nodes towards regions of the embedding space where they are maximally discriminative with respect to their contributions to positive and negative edges. To this end we take additional inspiration from triplet ranking loss functions [48] that are explicitly designed for learning representations that can capture relative similarities or differences between items.

With these considerations in mind, we initially posit the energy function  $\ell_{node} \triangleq \|Y - f(X; W)\|_F^2 +$

$$\lambda \sum_{(i,j) \in \mathcal{E}} \left[ \text{dist}(y_i, y_j) - \frac{\lambda_K}{\lambda K} \sum_{j' \in \mathcal{V}_{(i,j)}^K} \text{dist}(y_i, y_{j'}) \right], \quad (2)$$

where  $f(X; W)$  (assumed to apply row-wise to each individual node feature  $x_i$ ) represents a base model that processes the input features using trainable weights  $W$ ,  $\text{dist}(y_i, y_j)$  is a distance metric, while  $\lambda$  and  $\lambda_K$  are hyperparameters that control the impact of positive and negative edges. Moreover,  $\mathcal{V}_{(i,j)}^K$  is the set of negative destination nodes sampled for edge  $(v_i, v_j)$  and  $|\mathcal{V}_{(i,j)}^K| = K$ . Overall, the first term pushes the embeddings towards the processed input features, while the second and third terms apply penalties to positive and negative edges in a way that is loosely related to the aforementioned triplet ranking loss (more on this below).

If we choose  $\text{dist}(i, j) = \|y_i - y_j\|^2$  and define edges of the negative graph  $\mathcal{G}^-$  as  $\mathcal{E}^- \triangleq \{(i, j') \mid j' \in \mathcal{V}_{(i,j)}^K \text{ for } (i, j) \in \mathcal{E}\}$ , we can rewrite (2) as  $\ell_{node} =$

$$\|Y - f(X; W)\|_F^2 + \lambda \text{tr}[Y^\top LY] - \frac{\lambda_K}{K} \text{tr}[Y^\top L^- Y], \quad (3)$$

where  $L^-$  is the Laplacian matrix of  $\mathcal{G}^-$ . To find the minimum, we compute the gradients

$$\frac{\partial \ell_{node}}{\partial Y} = 2(Y - f(X; W)) + 2\lambda LY - \frac{\lambda_K}{K} 2L^- Y, \quad (4)$$

where  $Y^{(0)} = f(X; W)$ . The corresponding gradient descent updates then become  $Y^{(t+1)} =$

$$Y^{(t)} - \alpha((Y^{(t)} - f(X; W)) + \lambda LY^{(t)} - \frac{\lambda_K}{K} L^- Y^{(t)}), \quad (5)$$

where the step size is  $\frac{\alpha}{2}$ . We note however that (3) need not generally be convex or even lower bounded. Moreover, the gradients may be poorly conditioned for fast convergence depending on the Laplacian matrices involved. Hence we consider several refinements next to stabilize the learning process.

### B. Energy Function Refinements

**Lower-Bounding the Negative Graph.** Since the regularization of negative edges brings the possibility of an ill-posed loss surface (a non-convex loss surface that may be unbounded from below), we introduce a convenient graph-aware lower-bound analogous to the max operator used by the triplet loss. Specifically, we update (3) to the form  $\ell_{node} = \|Y - f(X; W)\|_F^2 +$

$$\lambda \text{tr}[Y^\top LY] + \frac{\lambda_K}{K} \text{Softplus}(|\mathcal{E}| \gamma - \text{tr}[Y^\top L^- Y]), \quad (6)$$

noting that we use  $\text{Softplus}(x) = \log(1 + e^x)$  instead of  $\max(\cdot, 0)$  to make the energy differentiable. Unlike triplet loss that includes positive term in the  $\max(\cdot, 0)$  function, we only restrict the lower bound for the negative term, because we still want the positive part to impact our model when the negative part hits the bound.

**Normalization.** We use the normalized Laplacian matrices of original graph  $\mathcal{G}$  and negative graph  $\mathcal{G}^-$  instead to make the gradients smaller. Also, for the gradients of the first term in (3), we apply a re-scaling by summation of degrees of both graphs. The modified gradients are  $\frac{\partial \ell_{node}}{\partial Y} =$

$$2(D + D^-)^{-1}(Y - f(X; W)) + 2\lambda \tilde{L} Y - \frac{\lambda_K}{K} 2\tilde{L}^- Y, \quad (7)$$

where  $D$  and  $D^-$  are diagonal degree matrices of  $\mathcal{G}$  and  $\mathcal{G}^-$ . The normalized Laplacians are  $\tilde{L} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$ ,  $\tilde{L}^- = D^{-\frac{1}{2}} L^- D^{-\frac{1}{2}}$  leading to the corresponding energy

$$\ell_{node} = \|(D + D^-)^{-1}(Y - f(X; W))\|_F^2 + \lambda \text{tr}[Y^\top \tilde{L} Y] - \frac{\lambda_K}{K} \text{tr}[Y^\top \tilde{L}^- Y]. \quad (8)$$

**Learning to Combine Negative Graphs.** We now consider a more flexible implementation of negative graphs. More concretely, we sample  $K$  negative graphs  $\{\mathcal{G}_{(k)}^-\}_{k=1, \dots, K}$ , in which every negative graph consists of one negative edge per positive edge ( $\mathcal{G}_{(k)}^- \triangleq \{(i, j') \mid j' \in \mathcal{V}_{(i,j)}^1 \text{ for } (i, j) \in \mathcal{E}\}$ ) (note that the superscript on  $\mathcal{V}$  implies a single negative sample per positive edge). And we set learnable weights  $\lambda_k^k$  for the

structure term of each negative graph, which converts the energy function to  $\ell_{node} =$

$$\|Y - f(X; W)\|_F^2 + \lambda \text{tr}[Y^\top LY] - \frac{1}{K} \sum_{k=1}^K \lambda_K^k \text{tr}[Y^\top L_k^- Y], \quad (9)$$

Also in practice we normalize this energy function to

$$\ell_{node} = \|(D + D_K^-)^{-1}(Y - f(X; W))\|_F^2 + \lambda \text{tr}[Y^\top \tilde{L} Y] - \frac{1}{K} \sum_{k=1}^K \lambda_K^k \text{tr}[Y^\top \tilde{L}_k^- Y], \quad (10)$$

where  $D_K^- = \sum_{k=1}^K D_k^-$ ,  $D_k^-$  is the degree matrix of  $L_k^-$  and  $\tilde{L}_k^- = D_K^{-\frac{1}{2}} L_k^- D_K^{-\frac{1}{2}}$ . The lower bound is also added as before. Overall, the motivation here is to inject trainable flexibility into the negative sample graph, which is useful for increasing model expressiveness.

### C. The Overall Algorithm

Combining the modifications we discussed in last section (and assuming a single, fixed  $\lambda_K$  here for simplicity; the more general, learnable case with multiple  $\lambda_K^k$  from Section IV-B naturally follows), we obtain the final energy function

$$\ell_{node} = \|(D + D^-)^{-1}(Y - f(X; W))\|_F^2 + \lambda \text{tr}[Y^\top \tilde{L} Y] + \frac{\lambda_K}{K} \text{Softplus}(\gamma|\mathcal{E}| - \text{tr}[Y^\top \tilde{L}^- Y]), \quad (11)$$

with the associated gradients  $\frac{\partial \ell_{node}}{\partial Y} =$

$$2(D + D^-)^{-1}(Y - f(X; W)) + 2\lambda \tilde{L} Y - \frac{\lambda_K}{K} 2\tilde{L}^- Y \sigma(Q), \quad (12)$$

where  $Q = \gamma|\mathcal{E}| - \text{tr}[Y^\top \tilde{L}^- Y]$  and  $\sigma(x)$  is the sigmoid function. The final updates for our model then become

$$\begin{aligned} Y^{(t+1)} &= Y^{(t)} - \alpha \left( (D + D^-)^{-1}(Y^{(t)} - f(X; W)) + \lambda \tilde{L} Y^{(t)} \right. \\ &\quad \left. - \frac{\lambda_K}{K} \tilde{L}^- Y^{(t)} \sigma(Q^{(t)}) \right) \\ &= C_1 Y^{(t)} + C_2 f(X; W) + c_3 \tilde{A} Y^{(t)} - c_4 \tilde{A}^- Y^{(t)}, \end{aligned} \quad (13)$$

where the diagonal scaling matrices ( $C_1, C_2$ ) and scalar coefficients ( $c_3, c_4$ ) are given by

$$\begin{aligned} C_1 &= \left(1 - \alpha\lambda + \frac{1}{K}\alpha\lambda_K\sigma(Q^{(t)})\right) I - \alpha(D + D^-)^{-1}, \\ C_2 &= \alpha(D + D^-)^{-1}, \\ c_3 &= \alpha\lambda, \quad c_4 = \left(\frac{1}{K}\alpha\lambda_K\sigma(Q^{(t)})\right), \end{aligned} \quad (14)$$

with  $Q^{(t)} = \gamma|\mathcal{E}| - \text{tr}[(Y^{(t)})^\top \tilde{L}^- Y^{(t)}]$ ,  $I$  as an  $n \times n$  identity matrix, and  $\frac{\alpha}{2}$  as the step size.

From the above expressions, we observe that the first and second terms of (13) can be viewed as rescaled skip connections from the previous layer and input layer/base model, respectively. As we will later show, these scale factors are designed to facilitate guaranteed descent of the objective from (11). Meanwhile, the third term of (13) represents a typical GNN graph propagation layer while the fourth term is the

analogous negative sampling propagation unique to our model. In the context of Chinese philosophy, the latter can be viewed as the Yin to the Yang which is the third term, and with a trade-off parameter that can be learned when training the higher-level objective from (1), the Yin/Yang balance can in a loose sense be estimated from the data; hence the name *YinYangGNN* for our proposed approach.

We illustrate key aspects of YinYangGNN in Figure 1 and the explicit computational steps are shown in Algorithm 1. And we emphasize that although the derivations of YinYangGNN may be somewhat complex, *the algorithm that results is quite simple and efficient.*

**YinYangGNN Training.** Upon inspection of Algorithm 1, we observe that lines 3 to 9 depict the training process for each epoch. Here we first sample the negative graph on line 3 which is used (unique to our model) on lines 4 to 6 to produce the YinYangGNN encoder node embeddings  $Y^{(T)}$ ; this completes one model *forward pass*. Next, on line 7 we again sample negative edges as needed by the supervised training loss. Given these negative samples, along with predictive probabilities obtained via node embeddings from  $Y^{(T)}$  passed through the HadamardMLP decoder, we compute the link prediction training loss defined by (1) on line 8. Finally, we update model parameters  $\{W, \theta\}$  (and optionally  $\lambda_K$ ) by standard backpropagation on line 9; this executes one model *backward pass* (in our experiments we use the Adam optimizer).

**YinYangGNN Inference.** For inference, we basically follow one forward pass of the training process from above to obtain encoder node embeddings  $Y^{(T)}$ . These are then applied to the HadamardMLP decoder to produce link prediction scores or edge probability estimates analogous to line 8, but without the loss computation.

---

#### Algorithm 1: YinYangGNN Bilevel Optimization Algorithm for Link Prediction

---

**Input:** Graph  $\mathcal{G}$ , node features  $X$ , number of layers  $T$ , number of epochs  $M$ , trainable base model  $f(\cdot; W)$ , HadamardMLP predictor  $h(\cdot; \theta)$   
**Output:** Trained model weights  $W$  and HadamardMLP predictor weights  $\theta$ .

- 1 Set initial prediction  $Y^{(0)} = f(X; W)$ , where  $f$  is the trainable base model.
- 2 **for**  $Epoch = 1$  to  $M$  **do**
- 3     Sample negative graph  $\mathcal{G}^-$ .
- 4     **for**  $t = 0$  to  $T - 1$  **do**
- 5          $Y^{(t+1)} = \text{Update}(Y^{(t)})$ , computed via (13).
- 6     **end**
- 7     Sample negative edges for supervision based on  $\mathcal{E}_{tr}$ .
- 8     Compute the link prediction loss (1) with node embeddings  $Y^{(T)}$  and HadamardMLP  $h(\cdot; \theta)$ .
- 9     Backpropagate over parameters  $\{W, \theta\}$  (and optionally  $\lambda_K$ ) using optimizer (Adam, etc.)
- 10 **end**

---

## V. ANALYSIS

In this section we first address the computational complexity and convergence issues of our model before turning to further

TABLE II

TIME COMPLEXITY COMPARISONS. FOR BUDDY,  $b$  IS THE HOP NUMBER FOR PROPAGATION AND  $h$  IS THE COMPLEXITY OF HASH OPERATIONS. \* INDICATES THAT THE COMPLEXITY CAN BE SUBLINEAR TO  $n$  VIA [16], AN OPTION THAT IS ONLY AVAILABLE TO NODE-WISE EMBEDDING MODELS SUCH AS YINYANGNN.

|            | SEAL                                      | BUDDY                                | YinYanGNN   | GCN   |
|------------|---|--------------------------------------|---|---|
| Preprocess | $O(1)$                                    | $O(b \mathcal{E} (d+h))$             | $O(1)$  | $O(1)$  |
| Train      | $O( \mathcal{E} d^2 \mathcal{E}_{tr} )$   | $O((b^2h+bd^2) \mathcal{E}_{tr} )$   | $O(T \mathcal{E} Kd+nPd^2+ \mathcal{E}_{tr} d^2)$ | $O(T( \mathcal{E} d+nd^2)+ \mathcal{E}_{tr} d^2)$ |
| Encode     | $O( \mathcal{E} d^2 \mathcal{V}_{src} n)$ | $O((b^2h+bd^2) \mathcal{V}_{src} n)$ | $O(T \mathcal{E} Kd+nPd^2)$                       | $O(T( \mathcal{E} d+nd^2))$                       |
| Decode     |   |                                      | $O(d^2 \mathcal{V}_{src} n)^*$                    | $O(d^2 \mathcal{V}_{src} n)^*$                    |

insights into the role of negative sampling in our proposed energy function.

### A. Time Complexity

YinYanGNN has a time complexity given by  $O(T|\mathcal{E}|Kd+nPd^2)$  for one forward pass, where as before  $n$  is the number of nodes,  $T$  is the number of propagation layers/iterations,  $P$  is the number of MLP layers in the base model  $f(X;W)$ , and  $d$  is the hidden embedding size. Notably, this complexity is of the same order as a vanilla GCN model, one of the most common GNN architectures [7], which has a forward-pass complexity of  $O(T(|\mathcal{E}|d+nd^2))$ .

We now drill down further into the details of overall inference speed. We denote the set of source nodes for test-time link prediction as  $\mathcal{V}_{src}$ , and for each node we examine all the other nodes in the graph, which means we have to calculate roughly  $|\mathcal{V}_{src}|n$  edges. We compare YinYanGNN’s time with SEAL [9] and a recently proposed fast baseline BUDDY [12] in Table II. Like other node-wise methods, we split the inference time of our model into two parts: computing embeddings and decoding (for SEAL and BUDDY they are implicitly combined). The node embedding computation can be done only once so it does not depend on  $\mathcal{V}_{src}$ . Our decoding process uses HadamardMLP to compute scores for each destination node (which can also be viewed as being for each edge) and get the top nodes (edges). From this it is straightforward to see that the decoding time dominates the overall computation of embeddings. So for the combined inference time, SEAL is the slowest because of the factor  $|\mathcal{E}|n$  while BUDDY and our model are both linear in the graph node number  $n$  independently of  $|\mathcal{E}|$ . However, BUDDY has a larger factor including the computation of subgraph structure  $b^2h$ , which will be much slower as our experiments will show. Moreover, unlike SEAL or BUDDY, we can apply Flashlight [16] to node-wise methods like YinYanGNN, an accelerated decoding method based on maximum inner product search (MIPS) that allows HadamardMLP to have sublinear decoding complexity in  $n$ . See Section VI-C for related experiments and discussion.

### B. Convergence of YinYanGNN Layers/Iterations

Convergence criteria for the energies from (3) and (11) are as follows (see Appendix X-B for proofs):

**Proposition V.1.** *If  $\lambda_K < K \cdot \delta_{max}$ , where  $\delta_{max}$  is the largest eigenvalue of  $L^-$ , then (3) has a unique global*

*minimum. Moreover, if the step-size parameter satisfies  $\alpha < \left\| I + \lambda \tilde{L} - \frac{\lambda_K}{K} \tilde{L}^- \right\|_F^{-1}$ , then the gradient descent iterations of (5) are guaranteed to converge to this minimum.*

**Proposition V.2.** *There exists an  $\alpha' > 0$  such that for any  $\alpha \in (0, \alpha']$ , the iterations (13) will converge to a stationary point of (11).*

### C. Role of Negative Sampling in Proposed Energy

Previous work [49], [50] has demonstrated how randomly dropping edges of two isomorphic graphs can produce non-isomorphic, distinguishable subgraphs, and in so doing, enhance GNN expressiveness. Our incorporation of negative samples within the YinYanGNN encoder/forward pass can loosely be viewed in similar terms, but with a novel twist: Instead of randomly *dropping* edges in the original graph, we randomly *add* typed negative edges that are likewise capable of breaking otherwise indistinguishable symmetries.

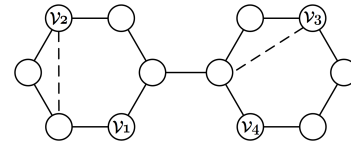


Fig. 2. Modified from [10]. Solid lines represent the original edges and dashed lines represent negative edges sampled in our model architecture (for simplicity we do not draw all negative edges).

Figure 2 serves to illustrate how the inclusion of negative samples within the forward pass of our model can potentially increase the expressiveness beyond traditional node-wise embedding approaches. As observed in the figure,  $v_2$  and  $v_3$  are isomorphic nodes in the original graph (solid lines). However, when negative samples/edges are included the isomorphism no longer holds, meaning that link  $(v_1, v_2)$  and link  $(v_1, v_3)$  can be distinguished by a node-wise embedding method even without unique discriminating input features. Moreover, when combined with the flexibility of learning to balance multiple negative sampling graphs as in (10) through the trainable weights  $\{\lambda_K^k\}$ , the expressiveness of YinYanGNN becomes strictly greater than or equal to a vanilla nodewise embedding method (with equivalent capacity) that has no explicit access to the potential symmetry-breaking influence of negative samples.

Critically though, these negative samples are not arbitrarily inserted into our modeling framework. Rather, they *emerge*

TABLE III

RESULTS ON LINK PREDICTION BENCHMARKS. BASELINE RESULTS ARE CITED FROM PRIOR WORK [12], [18] AND THE OGB LEADERBOARD (COMPARISONS WITH ADDITIONAL BASELINES CAN BE FOUND IN THE APPENDIX X-A). "-" MEANS NOT REPORTED. THE FORMAT IS AVERAGE SCORE  $\pm$  STANDARD DEVIATION. THE BEST RESULTS ARE BOLD-FACED AND UNDERLINED. OOM MEANS OUT OF GPU MEMORY.

|                   |                  | Cora                             | Citeseer                         | Pubmed                           | Collab                           | PPA                              | Citation2                        | DDI                              |
|-------------------|------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
|                   |                  | HR@100                           | HR@100                           | HR@100                           | HR@50                            | HR@100                           | MRR                              | HR@20                            |
| Edge-wise<br>GNNs | <b>SEAL</b>      | 81.71 $\pm$ 1.30                 | 83.89 $\pm$ 2.15                 | 75.54 $\pm$ 1.32                 | 64.74 $\pm$ 0.43                 | 48.80 $\pm$ 3.16                 | 87.67 $\pm$ 0.32                 | 30.56 $\pm$ 3.86                 |
|                   | <b>NBFnet</b>    | 71.65 $\pm$ 2.27                 | 74.07 $\pm$ 1.75                 | 58.73 $\pm$ 1.99                 | OOM                              | OOM                              | OOM                              | 4.00 $\pm$ 0.58                  |
|                   | <b>Neo-GNN</b>   | 80.42 $\pm$ 1.31                 | 84.67 $\pm$ 2.16                 | 73.93 $\pm$ 1.19                 | 57.52 $\pm$ 0.37                 | 49.13 $\pm$ 0.60                 | 87.26 $\pm$ 0.84                 | 63.57 $\pm$ 3.52                 |
|                   | <b>GDGNN</b>     | —                                | —                                | —                                | 54.74 $\pm$ 0.48                 | 45.92 $\pm$ 2.14                 | 86.96 $\pm$ 0.28                 | —                                |
|                   | <b>SUREL</b>     | —                                | —                                | —                                | 63.34 $\pm$ 0.52                 | 53.23 $\pm$ 1.03                 | <b>89.74<math>\pm</math>0.18</b> | —                                |
|                   | <b>SUREL+</b>    | —                                | —                                | —                                | 64.10 $\pm$ 1.06                 | 54.32 $\pm$ 0.44                 | 88.90 $\pm$ 0.06                 | —                                |
|                   | <b>BUDDY</b>     | 88.00 $\pm$ 0.44                 | 92.93 $\pm$ 0.27                 | 74.10 $\pm$ 0.78                 | 65.94 $\pm$ 0.58                 | 49.85 $\pm$ 0.20                 | 87.56 $\pm$ 0.11                 | 78.51 $\pm$ 1.36                 |
| Node-wise<br>GNNs | <b>GCN</b>       | 66.79 $\pm$ 1.65                 | 67.08 $\pm$ 2.94                 | 53.02 $\pm$ 1.39                 | 44.75 $\pm$ 1.07                 | 18.67 $\pm$ 1.32                 | 84.74 $\pm$ 0.21                 | 37.07 $\pm$ 5.07                 |
|                   | <b>SAGE</b>      | 55.02 $\pm$ 4.03                 | 57.01 $\pm$ 3.74                 | 39.66 $\pm$ 0.72                 | 48.10 $\pm$ 0.81                 | 16.55 $\pm$ 2.40                 | 82.60 $\pm$ 0.36                 | 53.90 $\pm$ 4.74                 |
|                   | <b>YinYanGNN</b> | <b>93.83<math>\pm</math>0.78</b> | <b>94.45<math>\pm</math>0.53</b> | <b>90.73<math>\pm</math>0.40</b> | <b>66.10<math>\pm</math>0.20</b> | <b>54.64<math>\pm</math>0.49</b> | 86.21 $\pm$ 0.09                 | <b>80.92<math>\pm</math>3.35</b> |

by taking gradient steps (12) over a principled regularization factor (within (10)) designed to push the embeddings of nodes sharing a negative edge apart during the forward pass. In a Section VII-E ablation we compare this unique YinYanGNN aspect with the alternative strategy of using random node features for breaking isomorphisms.

We conclude here by remarking that the goal of YinYanGNN, and the supporting analysis of this section, is not to explicitly match the expressiveness of edge-wise link prediction methods; even in principle this is not feasible given the additional information edge-wise methods have access to. Instead, the guiding principle of YinYanGNN is more aptly distilled as follows: conditioned on the limitations of a purely node-wise link prediction architecture, to what extent can we improve model expressiveness (and ultimately accuracy) without compromising the scalability which makes node-wise methods attractive in the first place.

## VI. EXPERIMENTS

### A. Experimental Setup

a) *Datasets and Evaluation Metrics.*: We evaluate YinYanGNN for link prediction on Planetoid datasets: Cora [51], Citeseer [52], Pubmed [53], and Open Graph Benchmark (OGB) link prediction datasets [45]: ogbl-collab, ogbl-PPA, ogbl-Citation2 and ogbl-DDI. Planetoid represents classic citation network data, whereas OGB involves challenging, multi-domain, diverse benchmarks involving large graphs. Detailed statistics are summarized in the Appendix X-C. We adopt the hits ratio @k(HR@k) as the main evaluation metric as in [12] for Planetoid datasets. This metric computes the ratio of positive edges ranked equal or above the  $k$ -th place out of candidate negative edges at test time. We set  $k$  to 100 for these three datasets. For OGB datasets, we follow the official settings. Note that the metric for ogbl-Citation2 is Mean Reciprocal Rank (MRR), meaning the reciprocal rank of positive edges among all the negative edges averaged over all source nodes. Finally, we choose the test results based on the best validation results. We also randomly select 10 different seeds and report

average results and standard deviation for all datasets. For implementation details and hyperparameters, please refer to Appendix X-C.

b) *Baseline Models.*: To calibrate the effectiveness of our model, in this section we conduct comprehensive comparisons with node-wise GNNs: GCN [7] and GraphSage [8], and edge-wise GNNs: SEAL [9], NeoGNN [17], NBFnet [11], BUDDY [12]), GDGNN [21], SUREL [22], and SUREL+ [18]. We differ to Appendix X-A additional experiments spanning traditional link prediction heuristic methods: Common Neighbors (CN) [4], Adamic-Adar (AA) [5] and Resource Allocation (RA) [6], non-GNN or graph methods: MLP, Node2vec [54], and Matrix-Factorization (MF) [1]), knowledge graph (KG) methods: TransE [55], ComplEx [56], and DistMult [57], additional node-wise GNNs: GAT [20], GIN [33], JKNet [58], and GCNII [59], and finally distillation methods.

Overall, for more standardized comparisons, we have chosen baselines based on published papers with open-source code and exclude those methods relying on heuristic augmentation strategies like anchor distances, or non-standard losses for optimization that can be applied more generally. Section VI-D also addresses additional link prediction baseline methods from the OGB leaderboard, including complementary methods that can be applied to YinYanGNN to further improve performance with additional complexity.

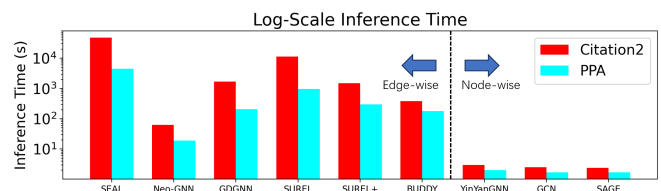


Fig. 3. Log-scale inference time. Citation2, PPA are the two largest OGB link prediction graphs.

TABLE IV  
DECODING SPEED AT INFERENCE TIME ON CITATION2 AS MEASURED BY NODES/SEC. (HIGHER IS BETTER).

|              | SEAL    | BUDDY | YinYanGNN | YinYanGNN(dot) | YinYanGNN(dot MIPS) |
|--------------|---------|-------|-----------|----------------|---------------------|
| nodes/second | 0.00024 | 0.005 | 0.6       | 1.3            | 25                  |

### B. Link Prediction Results

Accuracy results are displayed in Table III, where we observe that YinYanGNN achieves the best performance on 6 out of 7 datasets (while remaining competitive across all 7) even when compared against more time-consuming or inference-inefficient edge-wise methods. And we outperform node-wise methods by a large margin (including several others shown in Appendix X-A), demonstrating that YinYanGNN can achieve outstanding predictive accuracy without sacrificing efficiency. Similarly, as included in the Appendix X-A, YinYanGNN also outperforms a variety of non-GNN link prediction baselines.

### C. Efficiency Analysis

**Inference Time.** We next present comparisons in terms of inference speed, which is often the key factor determining whether or not a model can be deployed in real-world scenarios. For example, in an online system, providing real-time recommendations may require quickly evaluating a large number of candidate links. Figure 3 reports the results, again relative to both edge- and node-wise baseline models. Noting the log-scale time axis, from these results we observe that YinYanGNN is significantly faster than all the edge-wise models, and nearly identical to the fellow node-wise approaches as expected. And for the fastest edge-wise model, Neo-GNN, YinYanGNN is still simultaneously more efficient (Figure 3) and also much more accurate (Table III). Additional related details and inference-time results are deferred to Appendix X-A. We also remark that, as a node-wise model, the efficiency of YinYanGNN can be further improved to sublinear complexity using Flashlight [16]; however, at the time of submission public Flashlight code was not available, so we defer consideration to future work.

**Decoding Time.** Besides full inference times, in practice, we can save the embeddings output from the encoder and only compare the decoding step, which is a key differentiator. To this end, we compare the decoding speed of YinYanGNN with SEAL, a highly-influential edge-wise model, and BUDDY, which represents a strong baseline with a good balance between accuracy and speed for edge-wise GNNs. We conduct the experiments according to [16]. Table IV displays the results, where our model achieves the fastest performance compared with these two baselines. We also note that although the time complexity of BUDDY is linear in  $n$  (which is the same as our model without Flashlight), the different scaling coefficients caused by the required subgraph information calculation still make BUDDY much slower than YinYanGNN.

Finally, although the code for Flashlight is not yet publicly available, we can mimic the acceleration it will produce as follows. As Flashlight is a repeated version of MIPS, which

can be applied to dot product prediction, we also run MIPS on our model with dot product as the decoder to see how fast MIPS can accelerate. We run the experiments on CPU and achieve a 20x speedup (compare last two columns of Table IV with and without MIPS).

**Training Time.** We note that as is often the case for link prediction, our focus is on efficient inference, since the number of candidate node pairs grows quadratically in the graph size. That being said, we nonetheless provide a representative comparison here of training times. Specifically, we measure the training times on OGBL-PPA using identical training sets and platforms. BUDDY requires approximately 3 hours, whereas our YinYanGNN model accomplishes the same training task in approximately 12.5 minutes, underscoring a substantial discrepancy in training efficiency.

### D. Comparisons with Recent OGB Leaderboard Models

**General Comments.** After the edge-wise SEAL algorithm was originally published, not many new scalable node-wise alternatives have been proposed, likely because it has been difficult to achieve competitive performance. For example, there are actually very few node-wise entries atop the OGB link prediction leaderboard; almost all are less efficient edge-wise methods or related.

In fact, over the course of preparing this submission, there was no model of any kind with OGB leaderboard performance above ours on all datasets, and *only a single published node-wise model on the leaderboard above ours on any dataset*: this is model CFLP w/ JKNet [60], and only when applied to DDI. But when we examine the CFLP paper, we find that their performance is much lower than YinYanGNN on other non-OGB datasets as shown in Table V.

TABLE V  
COMPARISON OF YINYANGNN WITH CFLP.

|                      | Cora              | Citeseer          | Pubmed            |
|----------------------|-------------------|-------------------|-------------------|
| <b>Model</b>         | HR@20             | HR@20             | HR@20             |
| <b>CFLP w/ JKNet</b> | 65.57±1.05        | 68.09±1.49        | 44.90±2.00        |
| <b>YinYanGNN</b>     | <b>84.10±3.23</b> | <b>90.52±0.72</b> | <b>72.38±5.31</b> |

**Neural Common Neighbors (NCN).** Another recent approach with strong leaderboard performance worth mentioning here is NCN [61], which relies on a node-wise encoder as with YinYanGNN. However, the key feature of NCN is a decoder based on collecting common neighbors of nodes spanning each candidate link, which leads to a significantly higher inference cost. Even so, the NCN decoder itself can be easily be integrated with any node-wise encoder, including



YinYanGNN, to potentially increase accuracy, albeit with a much higher inference time. Results on OGB-Citation2 (the largest OGB link prediction task) are listed in Table VI. These results demonstrate that YinYanGNN with an NCN decoder can actually outperform both the original NCN MRR from [61] as well as all of the edge-wise models in Table III. The main downside is that inference time increases by  $\sim 40\times$  (see Table VI) over YinYanGNN in its original form. This further highlights the challenge of achieving the highest accuracy while preserving the full efficiency of purely node-wise models.

TABLE VI  
BETTER PERFORMANCES ON CITATION2 WITH NCN

|                  | Inference Time (s) | Citation2 (MRR)   |
|------------------|--------------------|-------------------|
| NCN [61]         | 116                | 88.64±0.14        |
| YinYanGNN        | 3                  | 86.21±0.09        |
| YinYanGNN w/ NCN | 117                | <b>90.03±0.02</b> |

**PLNLP and GIDN.** Beyond the above, we conclude by commenting on two additional leaderboard methods. First, there is an additional, unpublished node-wise model on the leaderboard called PLNLP [42] that does lead to strong accuracy on both Collab and DDI (two relatively small datasets), but the main contribution is a new training loss that is orthogonal to our method. In fact any node-wise model, including ours, could be trained with their loss to improve performance. The only downside is that the PLNLP approach seems to be dataset dependent, with lesser performance on other benchmarks, e.g., PLNLP only achieves 32.38 Hits@100 on PPA and 84.92 MRR on Citation2. Additionally, the GIDN model [62] follows the same loss as PLNLP and again achieves good results restricted to the smaller-scale Collab and DDI benchmarks (no other results are shown). However, the companion GIDN paper [62] is unpublished and extremely brief, with just a couple short paragraphs explaining the method and no assessment of inference time complexity. Hence we are unsure of how this model compares against node-wise models when applied to large-scale datasets as is our focus.

VII. ABLATIONS

A. Ablation over Negative Sampling.

As the integration of both positive and negative samples within a unified node-wise embedding framework is a critical component of our model, in Table VII we report results both with and without the inclusion of the negative sampling penalty in our lower-level embedding model from (13). Clearly the former displays notably superior performance as expected.

B. Learning Negative Graphs

One benefit of learning  $\{\lambda_K^k\}$  is that it allows us to better match the performance of larger  $K$  with fixed  $\lambda_K$  using a smaller  $K$  but with learnable  $\{\lambda_K^k\}$ . The latter is more practical given the greater efficiency of a smaller  $K$ . Figures 4 and 5 demonstrate this phenomena, whereby a learnable  $\{\lambda_K^k\}$  achieves better accuracy for smaller values of  $K$  on Cora and Pubmed datasets (for larger  $K$  performance is similar).

TABLE VII  
PERFORMANCE OF YINYANGNN WITH OR WITHOUT NEGATIVE SAMPLING IN MODEL ARCHITECTURE.

|              | Pubmed            | Collab            | PPA               | Citation2         | DDI               |
|--------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| YinYanGNN    | HR@100            | HR@50             | HR@100            | MRR               | HR@20             |
| W/O Negative | 86.70±3.23        | 63.02±0.44        | 49.36±2.91        | 83.45±0.21        | 57.80±5.39        |
| W/ Negative  | <b>90.73±0.40</b> | <b>66.10±0.20</b> | <b>54.64±0.49</b> | <b>86.21±0.09</b> | <b>80.92±3.35</b> |

For bigger datasets where using larger  $K$  can be prohibitively expensive, this capability can be exploited to achieve higher accuracy. For example, on PPA and Citation2 a larger  $K$  may well improve accuracy, but this is not computationally cheap without sampling or other approximations. In contrast, we can improve performance with small  $K$  just by learning  $\{\lambda_K^k\}$  as shown in Table VIII. (In contrast, for dense graphs a larger  $K$  is less likely to be helpful, in part because the extra edges will make the negative graph much more dense and more similar to complete graph, which results in the loss of structural information. Not surprisingly then, we found that learning  $\{\lambda_K^k\}$  did not improve performance on the dense graph DDI dataset where  $K = 1$  was optimal.)

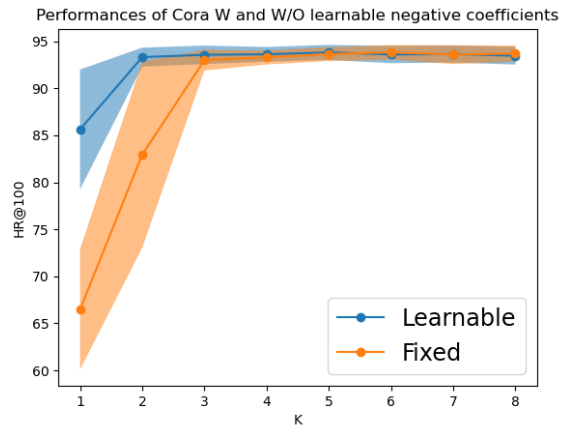


Fig. 4. Learnable or Fixed  $\{\lambda_K^k\}$  performances on Cora, standard deviation shown by backgrounds.

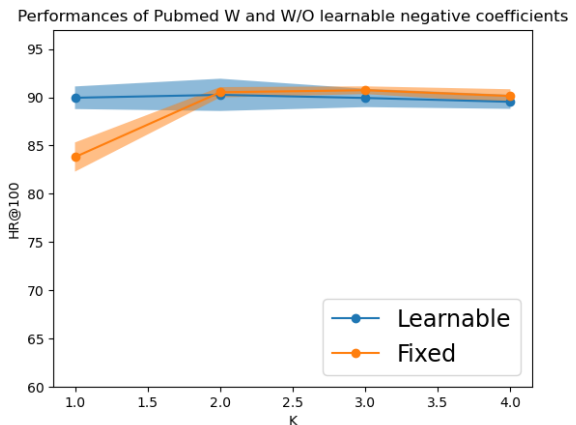


Fig. 5. Learnable or Fixed  $\{\lambda_K^k\}$  performances on Pubmed, standard deviation shown by backgrounds.

TABLE VIII  
BETTER PERFORMANCES ON PPA AND CITATION2 WITH LEARNABLE  $\{\lambda_K^k\}$

|                                    | PPA               | Citation2         |
|------------------------------------|-------------------|-------------------|
|                                    | HR@100            | MRR               |
| <b>fixed</b> $\{\lambda_K^k\}$     | 52.32±0.24        | 85.80±0.08        |
| <b>learnable</b> $\{\lambda_K^k\}$ | <b>54.64±0.49</b> | <b>86.21±0.09</b> |

C. Different Negative Samplers

As shown in Table IX, the YinYanGNN architecture is tolerant to the use of different negative samplers, from sampling negative edges uniformly in the graph (namely global uniform sampler) to sampling negative edges by fixing the source node and sampling negative destination nodes (namely source uniform sampler).

TABLE IX  
PERFORMANCE OF YINYANGNN WITH DIFFERENT SAMPLERS. GLOBAL MEANS GLOBAL UNIFORMLY SAMPLER. SOURCE MEANS SOURCE UNIFORMLY SAMPLER.

|         | Cora       | Citeseer   | Pubmed     |
|---------|------------|------------|------------|
| Sampler | HR@100     | HR@100     | HR@100     |
| Global  | 92.63±1.17 | 93.43±0.49 | 90.54±0.49 |
| Source  | 92.69±0.21 | 94.45±0.53 | 90.40±0.28 |

D. Negative sampling in the forward pass of common GNNs

As we discussed in Section IV, heuristically embedding negative samples within an existing GNN architecture is another possible strategy. But in this way the negative sampling terms are not integrated by a principled energy function. We follow the analysis of negative sampling in energy function, and modify the GCN architecture to incorporate negative sampling in the forward model as follows

$$Y^{(t+1)} = \text{ReLU} \left[ \left( \tilde{A} - \frac{\lambda_K}{K} \tilde{A}^- \right) Y^{(t)} W^{(t)} \right], \quad (15)$$

where  $W^{(t)}$  are the layer-wise weights and  $K$  is the negative edge number for each positive edge. Also,  $\lambda_K$  is the hyperparameter to control the impact of negative edges as before. As shown in Table X, negative sampling is also useful in some datasets for GCNs, but the performance is still not as good compared to YinYanGNN, where the negative samples are anchored to the proposed energy function.

E. Random Features versus Negative edges

In Section V-C, we show that random edges can in principle break the stated isomorphism. However, another potential way to accomplish this effect is to use random features that disambiguate each node. However, this can introduce undesirable auxiliary effects, which become apparent when we analyze the resulting energy function when applying random features instead of negative sampling.

TABLE X  
PERFORMANCE OF GCN WITH NEGATIVE EDGE SAMPLING, COMPARED WITH VANILLA GCN AND YINYANGNN.

|                 | Cora              | Citeseer          | Pubmed            |
|-----------------|-------------------|-------------------|-------------------|
| model           | HR@100            | HR@100            | HR@100            |
| GCN             | 66.79±1.65        | 67.08±2.94        | 53.02±1.39        |
| GCN W/ negative | 76.15±0.69        | 67.80±3.74        | 74.96±1.15        |
| YinYanGNN       | <b>93.83±0.78</b> | <b>94.45±0.53</b> | <b>90.73±0.40</b> |

When we add random features to (3) we obtain the modified form

$$\ell_{node_r}(Y) = \|Y - f(X + X_R; W)\|_F^2 + \lambda \text{tr}[Y^T L Y]. \quad (16)$$

For simplicity of illustration, we consider an  $f$  that obeys the distributive property, i.e., We then compute the gradient

$$\begin{aligned} \frac{\partial \ell_{node_r}}{\partial Y} &= 2(Y - f(X + X_R; W)) + 2\lambda L Y \\ &= 2(Y - f(X; W)) + 2\lambda L Y - 2f(X_R; W). \end{aligned} \quad (17)$$

From these expressions we observe that random features merely introduce random blurriness to the gradient updates; they do not emerge from an otherwise useful regularization factor. In contrast, with YinYanGNN, the negative samples create gradients that push the node embeddings of nodes that do not share a true edge apart, a useful form of structural regularization that is independent of any mechanism to break isomorphisms per se.

We support these points with additional experiments in Table XI, where random features are introduced and compared against our original YinYanGNN model. From these results we observe that random features can actually degrade performance and introduce instabilities as evidenced by the reduced accuracy and larger standard deviations.

TABLE XI  
PERFORMANCE OF YINYANGNN VARIANTS.

|                    | Cora              | Citeseer          | Pubmed            |
|--------------------|-------------------|-------------------|-------------------|
| YinYanGNN          | HR@100            | HR@100            | HR@100            |
| W/O Negative edges | 91.72±0.33        | 92.61±0.31        | 86.70±3.23        |
| W/ Random features | 91.25±1.41        | 89.69±0.78        | 79.82±5.15        |
| W/ Negative edges  | <b>93.83±0.78</b> | <b>94.45±0.53</b> | <b>90.73±0.40</b> |

F. Proportion of Negative Graphs

And finally, in Table XII, we examine how YinYanGNN performance is impacted as we vary the relative proportion assigned to negative graphs, which amounts to the relative weighting of  $\lambda$  to  $\lambda_K$ . Specifically, we choose  $\lambda = 1$  and vary  $\lambda_K$  from 0.1 to 10 (in each case  $\lambda_K$  is fixed, not learnable) and train YinYanGNN on small datasets Cora, Citeseer, and Pubmed as well as large dataset Citation2. Performance is generally best with  $\lambda_K \in [0.1, 1]$  across all datasets, with the smaller datasets favoring  $\lambda_K = 1$  and the largest  $\lambda_K = 0.1$ .

TABLE XII  
PERFORMANCE OF YINYANGNN WITH DIFFERENT NEGATIVE GRAPH PROPORTIONS.

|                       | <b>Cora</b>                      | <b>Citeseer</b>                  | <b>Pubmed</b>                    | <b>Citation2</b>                 |
|-----------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| $\lambda : \lambda_K$ | HR@100                           | HR@100                           | HR@100                           | MRR                              |
| 1:0.1                 | 93.26 $\pm$ 0.95                 | 92.44 $\pm$ 0.59                 | 87.08 $\pm$ 4.45                 | <b>85.80<math>\pm</math>0.08</b> |
| 1:1                   | <b>93.83<math>\pm</math>0.78</b> | <b>94.45<math>\pm</math>0.53</b> | <b>90.73<math>\pm</math>0.40</b> | 81.76 $\pm$ 0.08                 |
| 1:3                   | 92.23 $\pm$ 1.08                 | 92.42 $\pm$ 1.10                 | 63.22 $\pm$ 2.63                 | 67.65 $\pm$ 0.08                 |
| 1:5                   | 91.54 $\pm$ 1.01                 | 92.45 $\pm$ 0.63                 | 63.20 $\pm$ 2.20                 | 58.75 $\pm$ 0.03                 |
| 1:10                  | 87.31 $\pm$ 1.81                 | 92.50 $\pm$ 0.34                 | 57.00 $\pm$ 1.24                 | 49.79 $\pm$ 0.12                 |

## VIII. CONCLUSION

In conclusion, we have proposed the YinYanGNN link prediction model that achieves accuracy on par with far more expensive edge-wise models, but with the efficiency of relatively cheap node-wise alternatives. This competitive balance is accomplished using a novel node-wise architecture that incorporates informative negative samples/edges into the design of the model architecture itself to increase expressiveness, as opposed to merely using negative samples for computing a training signal as in prior work. Given the critical importance of inference speed in link prediction applications, YinYanGNN represents a promising candidate for practical usage. In terms of limitations, we have not integrated our implementation with Flashlight for maximally accelerated inference, as public code is not yet available.

## IX. ACKNOWLEDGMENTS

This work was supported by the National Key Research and Development Program of China (No.2022ZD0160102), and conducted while the first author was an intern at Amazon.

## REFERENCES

- [1] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [2] B. P. Chamberlain, E. Rossi, D. Shiebler, S. Sedhain, and M. M. Bronstein, "Tuning word2vec for large scale recommendation systems," in *RecSys*, 2020.
- [3] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. Van Den Berg, I. Titov, and M. Welling, "Modeling relational data with graph convolutional networks," in *The semantic web: 15th international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings 15*. Springer, 2018, pp. 593–607.
- [4] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *Science*, vol. 286, no. 5439, pp. 509–512, 1999.
- [5] L. A. Adamic and E. Adar, "Friends and neighbors on the web," *Social networks*, vol. 25, no. 3, pp. 211–230, 2003.
- [6] T. Zhou, L. Lü, and Y.-C. Zhang, "Predicting missing links via local information," *The European Physical Journal B*, vol. 71, no. 4, pp. 623–630, oct 2009.
- [7] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.
- [8] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NeurIPS*, 2017.
- [9] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," in *NeurIPS*, 2018.
- [10] M. Zhang, P. Li, Y. Xia, K. Wang, and L. Jin, "Labeling trick: A theory of using graph neural networks for multi-node representation learning," *NeurIPS*, vol. 34, pp. 9061–9073, 2021.
- [11] Z. Zhu, Z. Zhang, L.-P. Xhonneux, and J. Tang, "Neural bellman-ford networks: A general graph neural network framework for link prediction," in *NeurIPS*, 2021.
- [12] B. P. Chamberlain, S. Shirobokov, E. Rossi, F. Frasca, T. Markovich, N. Y. Hammerla, M. M. Bronstein, and M. Hansmire, "Graph neural networks for link prediction with subgraph sketching," in *ICLR*, 2022.
- [13] A. Shrivastava and P. Li, "Asymmetric LSH (ALSH) for Sublinear Time Maximum Inner Product Search (MIPS)," in *Advances in Neural Information Processing Systems*, vol. 27, 2014.
- [14] B. Neyshabur and N. Srebro, "On Symmetric and Asymmetric LSHs for Inner Product Search," in *Proceedings of the 32nd International Conference on Machine Learning*. PMLR, Jun. 2015, pp. 1926–1934.
- [15] H.-F. Yu, C.-J. Hsieh, Q. Lei, and I. S. Dhillon, "A Greedy Approach for Budgeted Maximum Inner Product Search," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [16] Y. Wang, B. Hooi, Y. Liu, T. Zhao, Z. Guo, and N. Shah, "Flashlight: Scalable link prediction with effective decoders," in *Learning on Graphs Conference*. PMLR, 2022, pp. 14–1.
- [17] S. Yun, S. Kim, J. Lee, J. Kang, and H. J. Kim, "Neo-gnns: Neighborhood overlap-aware graph neural networks for link prediction," in *NeurIPS*, 2021.
- [18] H. Yin, M. Zhang, J. Wang, and P. Li, "Surel+: Moving from walks to sets for scalable subgraph-based graph representation learning," *Proceedings of the VLDB Endowment*, vol. 16, no. 11, pp. 2939–2948, 2023.
- [19] Y. Yang, T. Liu, Y. Wang, J. Zhou, Q. Gan, Z. Wei, Z. Zhang, Z. Huang, and D. Wipf, "Graph neural networks inspired by classical iterative algorithms," in *International Conference on Machine Learning*, 2021.
- [20] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv:1710.10903*, 2017.
- [21] L. Kong, Y. Chen, and M. Zhang, "Geodesic graph neural network for efficient graph representation learning," in *NeurIPS*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35, 2022, pp. 5896–5909.
- [22] H. Yin, M. Zhang, Y. Wang, J. Wang, and P. Li, "Algorithm and system co-design for efficient subgraph-based graph representation learning," *Proceedings of the VLDB Endowment*, vol. 15, no. 11, pp. 2788–2796, 2022.
- [23] H. Ahn, Y. Yang, Q. Gan, T. Moon, and D. Wipf, "Descent steps of a relation-aware energy produce heterogeneous graph neural networks," *arXiv:2206.11081*, 2022.
- [24] S. Chen and Y. C. Eldar, "Graph signal denoising via unrolling networks," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5290–5294.
- [25] X. Liu, W. Jin, Y. Ma, Y. Li, H. Liu, Y. Wang, M. Yan, and J. Tang, "Elastic graph neural networks," in *International Conference on Machine Learning*, 2021.
- [26] Y. Ma, X. Liu, T. Zhao, Y. Liu, J. Tang, and N. Shah, "A unified view on graph neural networks as graph signal denoising," *arXiv:2010.01777*, 2020.
- [27] X. Pan, S. Song, and G. Huang, "A unified framework for convolution-based graph neural networks," <https://openreview.net/forum?id=zUMD-Fb9Bt>, 2021. [Online]. Available: <https://openreview.net/forum?id=zUMD-Fb9Bt>
- [28] H. Zhang, T. Yan, Z. Xie, Y. Xia, and Y. Zhang, "Revisiting graph convolutional network on semi-supervised node classification from an optimization perspective," *CoRR*, 2020.
- [29] M. Zhu, X. Wang, C. Shi, H. Ji, and P. Cui, "Interpreting and unifying graph neural networks with an optimization framework," *arXiv:2101.11859*, 2021.
- [30] Z. Wang, Q. Ling, and T. Huang, "Learning deep  $\ell_0$  encoders," in *AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [31] J. Feng, Y. Chen, F. Li, A. Sarkar, and M. Zhang, "How powerful are k-hop message passing graph neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 4776–4790, 2022.
- [32] F. Frasca, B. Bevilacqua, M. Bronstein, and H. Maron, "Understanding and extending subgraph gnns by rethinking their symmetries," *Advances in Neural Information Processing Systems*, vol. 35, pp. 31376–31390, 2022.
- [33] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" *arXiv:1810.00826*, 2018.
- [34] B. Zhang, S. Luo, L. Wang, and D. He, "Rethinking the expressive power of gnns via graph biconnectivity," in *The Eleventh International Conference on Learning Representations*, 2023.
- [35] B. Zhang, G. Feng, Y. Du, D. He, and L. Wang, "A complete expressiveness hierarchy for subgraph GNNs via subgraph weisfeiler-lehman tests," in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 23–29 Jul 2023, pp. 41019–41077. [Online]. Available: <https://proceedings.mlr.press/v202/zhang23k.html>

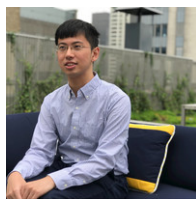
- [36] C. Morris, F. Geerts, J. Tönshoff, and M. Grohe, “WL meet VC,” in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 23–29 Jul 2023, pp. 25 275–25 302. [Online]. Available: <https://proceedings.mlr.press/v202/morris23a.html>
- [37] J. Zhou, J. Feng, X. Wang, and M. Zhang, “Distance-restricted folklore weisfeiler-leman gnns with provable cycle counting power,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [38] B. Srinivasan and B. Ribeiro, “On the equivalence between positional node embeddings and structural graph representations,” in *International Conference on Learning Representations*, 2020.
- [39] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen, “Graph contrastive learning with augmentations,” in *NeurIPS*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33, 2020, pp. 5812–5823.
- [40] J. Qiu, Q. Chen, Y. Dong, J. Zhang, H. Yang, M. Ding, K. Wang, and J. Tang, “Gcc: Graph contrastive coding for graph neural network pre-training,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ser. KDD ’20. New York, NY, USA: Association for Computing Machinery, 2020, p. 1150–1160.
- [41] W. Shiao, Z. Guo, T. Zhao, E. E. Papalexakis, Y. Liu, and N. Shah, “Link prediction with non-contrastive learning,” *arXiv:2211.14394*, 2022.
- [42] Z. Wang, Y. Zhou, L. Hong, Y. Zou, and H. Su, “Pairwise learning for neural link prediction,” *arXiv:2112.02936*, 2021.
- [43] S. Rendle, W. Krichene, L. Zhang, and J. Anderson, “Neural collaborative filtering vs. matrix factorization revisited,” in *Fourteenth ACM conference on recommender systems*, 2020, pp. 240–248.
- [44] C. Sun and G. Wu, “Adaptive graph diffusion networks with hop-wise attention,” *arXiv:2012.15024*, 2020.
- [45] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, “Open graph benchmark: Datasets for machine learning on graphs,” in *NeurIPS*, 2020.
- [46] J. Chen, J. Mueller, V. N. Ioannidis, S. Adeshina, Y. Wang, T. Goldstein, and D. Wipf, “Does your graph need a confidence boost? Convergent boosted smoothing on graphs with tabular node features,” in *International Conference on Learning Representations*, 2021.
- [47] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf, “Learning with local and global consistency,” *Advances in Neural Information Processing Systems*, 2004.
- [48] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, “Bpr: Bayesian personalized ranking from implicit feedback,” *arXiv:1205.2618*, 2012.
- [49] P. A. Papp, K. Martinkus, L. Faber, and R. Wattenhofer, “Dropgmn: Random dropouts increase the expressiveness of graph neural networks,” in *NeurIPS*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021, pp. 21 997–22 009.
- [50] F. Frasca, B. Bevilacqua, M. Bronstein, and H. Maron, “Understanding and extending subgraph gnns by rethinking their symmetries,” in *NeurIPS*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35, 2022, pp. 31 376–31 390.
- [51] A. K. McCallum, K. Nigam, J. Rennie, and K. Seymore, “Automating the construction of internet portals with machine learning,” *Information Retrieval*, vol. 3, no. 2, pp. 127–163, 2000.
- [52] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, “Collective classification in network data,” *AI Magazine*, vol. 29, no. 3, pp. 93–93, 2008.
- [53] G. Namata, B. London, L. Getoor, B. Huang, and U. EDU, “Query-driven active surveying for collective classification,” in *Proceedings Mining and Learning with Graphs*, 2012.
- [54] A. Grover and J. Leskovec, “node2vec: Scalable feature learning for networks,” in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 2016, pp. 855–864.
- [55] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, “Translating embeddings for modeling multi-relational data,” *NeurIPS*, vol. 26, 2013.
- [56] T. Trouillon, J. Welbl, S. Riedel, É. Gaussier, and G. Bouchard, “Complex embeddings for simple link prediction,” in *ICML*. PMLR, 2016, pp. 2071–2080.
- [57] B. Yang, W.-t. Yih, X. He, J. Gao, and L. Deng, “Embedding entities and relations for learning and inference in knowledge bases,” *arXiv:1412.6575*, 2014.
- [58] K. Xu, C. Li, Y. Tian, T. Sonobe, K.-i. Kawarabayashi, and S. Jegelka, “Representation learning on graphs with jumping knowledge networks,” in *ICML*. PMLR, 2018, pp. 5453–5462.
- [59] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li, “Simple and deep graph convolutional networks,” in *ICML*. PMLR, 2020, pp. 1725–1735.
- [60] T. Zhao, G. Liu, D. Wang, W. Yu, and M. Jiang, “Learning from counterfactual links for link prediction,” in *ICML*. PMLR, 2022, pp. 26 911–26 926.
- [61] X. Wang, H. Yang, and M. Zhang, “Neural common neighbor with completion for link prediction,” in *ICLR*, 2023.
- [62] Z. Wang, Y. Guo, J. Zhao, Y. Zhang, H. Yu, X. Liao, H. Jin, B. Wang, and T. Yu, “Gidn: A lightweight graph inception diffusion network for high-efficient link prediction,” *arXiv:2210.01301*, 2022.
- [63] Z. Guo, W. Shiao, S. Zhang, Y. Liu, N. V. Chawla, N. Shah, and T. Zhao, “Linkless link prediction via relational distillation,” in *ICML*. PMLR, 2023, pp. 12 012–12 033.
- [64] S. Zhang, Y. Liu, Y. Sun, and N. Shah, “Graph-less neural networks: Teaching old mlps new tricks via distillation,” in *ICLR*, 2021.
- [65] S. Bubeck *et al.*, “Convex optimization: Algorithms and complexity,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 3-4, pp. 231–357, 2015.
- [66] D. Bertsekas, *Nonlinear Programming*. Athena Scientific, 1999.



**Yuxin Wang** received the BS degree from Fudan University and is currently a Phd student in Fudan University major in computer science. His main research interest is efficiency in graph neural networks and natural language process.



**Xiannian Hu** received the MS degree in computer science from Fudan University in 2024. He is now working in TAOBAO & TMALL GROUP for pricing strategy and AI marketing, interest in machine learning, causal inference, operations optimization.



**Quan Gan** is a Senior Applied Scientist at Amazon. His research interests include heterogeneous graph neural networks, algorithm unrolling, and the application of graph neural networks to domains such as tabular data and hypergraph networks. He obtained his Master's degree at New York University and Bachelor's degree at Fudan University.



**Xuanjing Huang** is currently a Professor at the School of Computer Science, Fudan University, Shanghai, China. She obtained the PhD degree in Computer Science from Fudan University in 1998. Her research interests focus on natural language processing and information retrieval, with a particular emphasis on sentiment analysis, information extraction, pre-trained language models, and the robustness and interpretability of NLP.



**Xipeng Qiu** is a professor at the School of Computer Science, Fudan University. He received his B.S. and Ph.D. degrees from Fudan University. His research interests include natural language processing and deep learning.



**David Wipf** is a Principal Research Scientist at Amazon Web Services and an IEEE Fellow. Prior to industry positions at Amazon and Microsoft Research, he received a B.S. degree with highest honors from the University of Virginia, and M.S. and Ph.D. degrees from the University of California, San Diego as an NSF Fellow in Vision and Learning in Humans and Machines. He was later an NIH Postdoctoral Fellow at the University of California, San Francisco. His current research interests include deep generative models and graph neural networks. He is the recipient

of numerous fellowships and awards including the 2012 Signal Processing Society Best Paper Award.

X. APPENDIX

A. Additional Experimental Results

In this section we show further details of inference time results, and compare YinYanGNN with more diverse classical baselines, expressive node-wise GNNs, and distillation methods.

1) *Inference Time*: We show the inference time of different models on the two largest OGB-datasets in Table XIII.

TABLE XIII  
INFERENCE TIME ON TEST SET.

| Dataset   | SEAL      | Neo-GNN | GDGNN  | SUREL   | SUREL+ | BUDDY | YinYanGNN | GCN   | SAGE  |
|-----------|-----------|---------|--------|---------|--------|-------|-----------|-------|-------|
| PPA       | ≈ 4500 s  | 19 s    | 210 s  | 972 s   | 299 s  | 178 s | 2 s       | 1.7 s | 1.7 s |
| Citation2 | ≈ 48000 s | 63 s    | 1680 s | 11367 s | 1516 s | 385 s | 3 s       | 2.5 s | 2.4 s |

2) *Impact of the Number of Negative Samples on Performance and Inference Time*: Although detailed inference time comparisons using PPA and Citation2 were shown in Figure 3 and Table XIII, we further explore how the number of negative samples can affect both inference time as well as the corresponding performance on larger graphs. From the results shown in Table XIV we observe that YinYanGNN is robust to the number of negative samples for large datasets. Hence to prioritize efficiency  $K = 1$  is a reasonable selection.

TABLE XIV  
YINYANGNN PERFORMANCE AND INFERENCE TIMES VERSUS  $K$  ON PPA AND CITATION2.

| $K$ | PPA               | PPA                | Citation2         | Citation2          |
|-----|-------------------|--------------------|-------------------|--------------------|
|     | HR@100            | Inference Time (s) | MRR               | Inference Time (s) |
| 1   | 54.64±0.49        | 0.21               | <b>86.21±0.09</b> | 1.17               |
| 2   | <b>56.60±0.38</b> | 0.29               | 85.95±0.06        | 1.35               |
| 3   | 55.80±0.42        | 0.35               | 85.85±0.05        | 1.52               |

3) *Non-GNN Approaches*: Herein we show results for classical approaches reported from [12] and OGB leaderboards, with the exception that results for MLP and Node2vec for Cora, Citeseer and Pubmed are obtained by our own implementation in Appendix X-C. These new baselines include traditional heuristic methods (Common Neighbors (CN) [4], Adamic-Adar (AA) [5] and Resource Allocation (RA) [6]), non-GNN methods (MLP, Node2vec [54], Matrix-Factorization (MF) [1]), and knowledge graph (KG) methods(TransE [55], ComplEx [56], DistMult [57]). YinYanGNN performance exceeds that of these methods.

TABLE XV  
RESULTS ON LINK PREDICTION BENCHMARKS. BASELINE RESULTS ARE CITED FROM PREVIOUS WORKS [12], [18] AND THE OGB LEADERBOARD, WITH A FEW EXCEPTIONS MENTIONED IN THE TEXT.

|               |                  | Cora              | Citeseer          | Pubmed            | Collab            | PPA               | Citation2         | DDI               |
|---------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|               |                  | HR@100            | HR@100            | HR@100            | HR@50             | HR@100            | MRR               | HR@20             |
| Heuristics    | CN               | 33.92±0.46        | 29.79±0.90        | 23.13±0.15        | 56.44±0.00        | 27.65±0.00        | 51.47±0.00        | 17.73±0.00        |
|               | AA               | 39.85±1.34        | 35.19±1.33        | 27.38±0.11        | 64.35±0.00        | 32.45±0.00        | 51.89±0.00        | 18.61±0.00        |
|               | RA               | 41.07±0.48        | 33.56±0.17        | 27.03±0.35        | 64.00±0.00        | 49.33±0.00        | 51.98±0.00        | 27.60±0.00        |
| Non-GNN       | MLP              | 16.66±0.24        | 21.49±0.30        | 24.85±1.14        | 19.27±1.29        | 0.46±0.00         | 29.06±0.16        | NA                |
|               | Node2vec         | 46.87±0.97        | 50.76±2.32        | 54.40±11.10       | 48.88±0.54        | 22.26±0.88        | 61.41±0.11        | 23.26±2.09        |
|               | MF               | 37.37±0.00        | 48.96±0.02        | 13.79±11.97       | 38.86±0.29        | 32.29±0.94        | 51.86±4.43        | 13.68±4.75        |
| KG            | transE           | 67.40±1.60        | 60.19±1.15        | 36.67±0.99        | 29.40±1.15        | 22.69±0.49        | 76.44±0.18        | 6.65±0.20         |
|               | complEx          | 37.16±2.76        | 42.72±1.68        | 37.80±1.39        | 53.91±0.50        | 27.42±0.49        | 72.83±0.38        | 8.68±0.36         |
|               | DistMult         | 41.38±2.49        | 47.65±1.68        | 40.32±0.89        | 51.00±0.54        | 28.61±1.47        | 66.95±0.40        | 11.01±0.49        |
| Node-wise GNN | <b>YinYanGNN</b> | <b>93.83±0.78</b> | <b>94.45±0.53</b> | <b>90.73±0.40</b> | <b>66.10±0.20</b> | <b>54.64±0.49</b> | <b>86.21±0.09</b> | <b>80.92±3.35</b> |

4) *More node-wise GNNs*: We also include experiments with more expressive GNN architectures, namely GAT [20], GIN [33], JKNet [58], and GCNII [59], that all can be applied to arbitrary graph structure and have readily available public implementations. We tune these baselines according to our implementation in Appendix X-C. Results are shown in the table below, where we observe that YinYanGNN still maintains its strong advantage.

TABLE XVI  
COMPARISON OF YINYANGNN WITH ADDITIONAL NODE-WISE GNNS.

|                  | <b>Cora</b>       | <b>Citeseer</b>   | <b>Pubmed</b>     |
|------------------|-------------------|-------------------|-------------------|
| <b>Model</b>     | HR@100            | HR@100            | HR@100            |
| <b>GCNII</b>     | 48.11±1.86        | 44.23±1.55        | 49.72±2.39        |
| <b>GIN</b>       | 57.23±2.87        | 69.19±0.41        | 50.97±1.61        |
| <b>JK-Net</b>    | 67.86±2.32        | 54.04±3.73        | 62.73±4.53        |
| <b>GAT</b>       | 79.10±0.41        | 64.56±2.88        | 41.81±2.01        |
| <b>YinYanGNN</b> | <b>93.83±0.78</b> | <b>94.45±0.53</b> | <b>90.73±0.40</b> |

TABLE XVII  
RESULTS COMPARED WITH DISTILLATION METHODS.

|                                | <b>Cora</b>       | <b>Citeseer</b>   | <b>Pubmed</b>     | <b>Collab</b>     |
|--------------------------------|-------------------|-------------------|-------------------|-------------------|
| <b>Model</b>                   | HR@20             | HR@20             | HR@20             | HR@50             |
| <b>Logit Matching</b>          | 74.72±4.27        | 72.44±1.52        | 42.78±3.15        | 35.97±0.96        |
| <b>Representation Matching</b> | 75.75±1.51        | 65.19±5.54        | 44.44±2.40        | 36.86±0.45        |
| <b>LLP</b>                     | 78.82±1.74        | 77.32±2.42        | 57.33±2.42        | 49.10±0.57        |
| <b>YinYanGNN</b>               | <b>85.39±2.38</b> | <b>84.96±2.51</b> | <b>61.02±6.87</b> | <b>66.10±0.20</b> |

5) *Distillation Methods*: Another promising method for faster inference time is to distill GNN models. Such distillation methods are complementary to our work (which mainly focuses on GNN architecture design), and could in principle be applied in parallel using YinYanGNN as a teacher model. While we reserve such an effort to future work, for now we compare performance of existing distillation pipelines including LLP [63] and two baseline distillation models, logit matching [64] and representation matching. Results are in Table XVII, where the splits for Cora, Citeseer and Pubmed follow [63], while for Collab we use official settings. From the table we observe that YinYanGNN performance is considerably higher.

B. Proofs

1) *Proof for Proposition V.1*: We first restate Proposition V.1 here.

**Proposition X.1.** *If  $\lambda_K < K \cdot \delta_{max}$ , where  $\delta_{max}$  is the largest eigenvalue of  $L^-$ , then (3) has a unique global minimum. Moreover, if the step-size parameter satisfies  $\alpha < \left\| I + \lambda \tilde{L} - \frac{\lambda_K}{K} \tilde{L}^- \right\|_F^{-1}$ , then the gradient descent iterations of (5) are guaranteed to converge to this minimum.*

*Proof.* We first demonstrate conditions whereby (3) is strongly-convex, in which case it will have a unique global minimum. We can accomplish this by showing that the smallest eigenvalue of the corresponding Hessian Matrix is non-negative [65]. In our case, the smallest eigenvalue can be computed as  $1 + \lambda \delta_{min}^+ - \frac{\lambda_K}{K} \delta_{max}$ , where  $\delta_{min}^+ = 0$  is the smallest eigenvalue of  $L$ , and  $\delta_{max}$  is the largest eigenvalue of  $L^-$ , which is non-negative. So to guarantee strong convexity, we require that  $1 - \frac{\lambda_K}{K} \delta_{max} > 0$  or equivalently that  $\lambda_K < K \cdot \delta_{max}$ .

Proceeding further, if the gradient of a strongly-convex function is Lipschitz-continuous, with Lipschitz constant  $C$ , it follows that gradient descent with step-size less than  $C^{-1}$  will converge to the unique global minimum. In our setting, (3) has Lipschitz continuous gradients with Lipschitz constant that can be inferred via the bound

$$\begin{aligned} \left\| \frac{\partial \ell_{node}}{\partial Y_1} - \frac{\partial \ell_{node}}{\partial Y_2} \right\|_F^2 &= \left\| 2 \left[ Y_1 - Y_2 + \lambda \tilde{L}(Y_1 - Y_2) - \frac{\lambda_K}{K} \tilde{L}^-(Y_1 - Y_2) \right] \right\|_F^2 \\ &\leq \left\| 2 \left[ I + \lambda \tilde{L} - \frac{\lambda_K}{K} \tilde{L}^- \right] \right\|_F^2 \|Y_1 - Y_2\|_F^2, \end{aligned} \tag{18}$$

which holds for any pair of points  $\{Y_1, Y_2\}$ . From this it follows that  $C = \left\| 2 \left[ I + \lambda \tilde{L} - \frac{\lambda_K}{K} \tilde{L}^- \right] \right\|_F$ , and hence, if our assumed step-size  $\frac{\alpha}{2}$  is chosen such that  $\alpha < \left\| I + \lambda \tilde{L} - \frac{\lambda_K}{K} \tilde{L}^- \right\|_F^{-1}$ , the result is established. □

2) *Proof for Proposition 5.2*: We first restate the Proposition 5.2 here before the proof.

**Proposition X.2.** *There exists an  $\alpha' > 0$  such that for any  $\alpha \in (0, \alpha']$ , the iterations (13) will converge to a stationary point of (11).*

*Proof.* We first show that the energy from (11) has a Lipschitz continuous gradient within the compact set  $\|Y\|_F^2 \leq b^2$ , where  $b > 0$  is some constant and the associated (local) Lipschitz constant denoted  $C_b$  will be a function of this  $b$ . To show this, we begin by calculating the secondary gradient for (11). For this purpose we first write the gradient in the element-wise form

$$\frac{\partial \ell_{node}}{\partial Y_{ij}} = 2(D + D^-)^{-1}(Y_{ij} - f(X; W)_{ij}) + 2\lambda \sum_k \tilde{L}_{ik} Y_{kj} - 2\frac{\lambda_K}{K} \sum_k \tilde{L}_{ik}^- Y_{kj} \sigma(Q). \quad (19)$$

So the secondary gradient is

$$\frac{\partial^2 \ell_{node}}{\partial Y_{ij}^2} = 2(D + D^-)^{-1} Y_{ij} + 2\lambda \tilde{L}_{ii} - 2\frac{\lambda_K}{K} \tilde{L}_{ii}^- \sigma(Q) - 2\frac{\lambda_K}{K} \tilde{L}_{ii}^- Y_{ij} \sigma(Q)(1 - \sigma(Q)) \frac{\partial Q}{\partial Y_{ij}}. \quad (20)$$

Recall that  $Q = \gamma|\mathcal{E}| - \text{tr}[Y^\top \tilde{L}^- Y]$ , so

$$\frac{\partial Q}{\partial Y_{ij}} = -2 \sum_k \tilde{L}_{ik}^- Y_{kj}. \quad (21)$$

Consequently the secondary gradient becomes

$$\begin{aligned} \frac{\partial^2 \ell_{node}}{\partial Y_{ij}^2} &= 2(D + D^-)^{-1} Y_{ij} + 2\lambda \tilde{L}_{ii} - 2\frac{\lambda_K}{K} \tilde{L}_{ii}^- \sigma(Q) - 2\frac{\lambda_K}{K} \tilde{L}_{ii}^- \sigma(Q)(1 - \sigma(Q)) \left(-2 \sum_k \tilde{L}_{ik}^- Y_{kj}\right) Y_{ij} \\ &\leq 2|(D + D^-)^{-1}|b + 2\lambda|\tilde{L}_{ii}| + |2\frac{\lambda_K}{K} \tilde{L}_{ii}^-| + 4\frac{\lambda_K}{K} |\tilde{L}_{ii}^-| \left(\sum_k |\tilde{L}_{ik}^-| b^2\right), \end{aligned} \quad (22)$$

where  $|\cdot|$  denotes the absolute value of scalar-valued quantities. The above inequality comes from the fact that for any  $\|Y\|_F^2 \leq b^2$ ,  $|Y_{ij}| \leq b$  and  $|\sum_k \tilde{L}_{ik}^- Y_{kj} Y_{ij}| \leq \sum_k |\tilde{L}_{ik}^- Y_{kj} Y_{ij}| \leq \sum_k |\tilde{L}_{ik}^-| b^2$ . Besides,  $\sigma(Q) < 1$  and  $1 - \sigma(Q) < 1$ . Since elements of the Hessian exist and are bounded, the corresponding gradients must be Lipschitz continuous within the stated constraint set.

Next, we note that for any initialization  $Y^{(0)}$  we can choose a  $b$  sufficiently large such that  $\ell_{node}(Y) > \ell_{node}(Y^{(0)})$  for all  $Y \in \mathcal{S}_b \triangleq \{Y' : Y' \in \mathbb{R}^{n \times d}, \|Y'\|_F^2 > b^2\}$ . To see this, note that for any  $Y$  with  $\|Y\|_F^2$  sufficiently large,  $\ell_{node}(Y)$  is unbounded from above given that the first two terms collectively are strongly convex, while the last term is bounded from below.

We then define  $\bar{\mathcal{S}}_b = \mathbb{R}^{n \times d} \setminus \mathcal{S}_b$  (i.e., the complement of  $\mathcal{S}_b$ ) for convenience. Therefore, given any initialization  $Y^{(0)}$ , there will exist a Lipschitz constant  $C_b$  such that any gradient descent iteration computed at points  $Y \in \bar{\mathcal{S}}_b$  using a step-size  $\frac{\alpha}{2}$  less than  $1/C_b$  (so  $\alpha' \triangleq 2/C_b$ ) is guaranteed to reduce  $\ell_{node}(Y)$ , and in doing so, the iteration will remain within  $\bar{\mathcal{S}}_b$ . Hence we can apply standard results [66], for the convergence of gradient descent applied to non-convex functions  $f$  to a stationary point provided that  $f$  is lower-bounded (as is (11)) and has Lipschitz-continuous gradients, and the step-size parameter is chosen to be less than the inverse of the corresponding Lipschitz constant. □

### C. Further Experimental Details

1) *Datasets Statistics:* Statistics of the datasets used in the main paper are presented in Table XVIII. The splits for Cora, Citeseer and Pubmed are 7:1:2 for training:validation:test, following [12]. Splits for OGB datasets are the official ones from [45].

TABLE XVIII  
DATASETS STATISTICS.

|             | <b>Cora</b> | <b>Citeseer</b> | <b>Pubmed</b> | <b>Collab</b> | <b>PPA</b> | <b>DDI</b> | <b>Citation2</b> |
|-------------|-------------|-----------------|---------------|---------------|------------|------------|------------------|
| Node number | 2,708       | 3,327           | 18,717        | 235,868       | 576,289    | 4,267      | 2,927,963        |
| Edge number | 5,278       | 4,676           | 44,327        | 1,285,465     | 30,326,273 | 1,334,889  | 30,561,187       |
| splits      | rand        | rand            | rand          | time          | throughput | time       | protein          |
| avg degree  | 3.9         | 2.74            | 4.5           | 5.45          | 52.62      | 312.84     | 10.44            |



2) *Implementation Hardware and Hyperparameters*: The experiments were generally conducted on RTX 4090(24GB) and A100(40GB) GPUs (A100 only for ogbl-Citation2). The inference time comparisons in Table XIII were conducted on an RTX A10G, while the implementation of Table XIV was done using an A100(40GB). Hyperparameters were selected using a grid search on Cora, Citeseer and Pubmed. On OGB datasets, we empirically try different hyperparameters based on the experiences from the small datasets. The searching grid is the same for our model and baseline experiments if needed. For these we follow consistent ranges with the code from prior work for YinYanGNN, i.e., learning rate (0.0001-0.01), hidden dimension (32-1024), dropout (0-0.8). For specified hyperparameters for YinYanGNN, the searching grid is K (1-4), propagation step (2-16),  $\alpha$  (0.01-1),  $\lambda_K$  (100-1 or Learnable),  $\lambda$  (100-1).