# Neural Graph Matching Network: Learning Lawler's Quadratic Assignment Problem with Extension to Hypergraph and Multiple-graph Matching

Runzhong Wang, *Student Member, IEEE,* Junchi Yan, *Senior Member, IEEE,*
and Xiaokang Yang, *Fellow, IEEE*

**Abstract**—Graph matching involves combinatorial optimization based on edge-to-edge affinity matrix, which can be generally formulated as Lawler's Quadratic Assignment Problem (QAP). This paper presents a QAP network directly learning with the affinity matrix (equivalently the association graph) whereby the matching problem is translated into a constrained vertex classification task. The association graph is learned by an embedding network for vertex classification, followed by Sinkhorn normalization and a cross-entropy loss for end-to-end learning. We further improve the embedding model on association graph by introducing Sinkhorn based matching-aware constraint, as well as dummy nodes to deal with unequal sizes of graphs. To our best knowledge, this is one of the first network to directly learn with the general Lawler's QAP. In contrast, recent deep matching methods focus on the learning of node/edge features in two graphs respectively. We also show how to extend our network to hypergraph matching, and matching of multiple graphs. Experimental results on both synthetic graphs and real-world images show its effectiveness. For pure QAP tasks on synthetic data and QAPLIB benchmark, our method can perform competitively and even surpass state-of-the-art graph matching and QAP solvers with notable less time cost. We provide a project homepage at http://thinklab.sjtu.edu.cn/project/NGM/index.html.

**Index Terms**—Graph Matching, Deep Learning, Quadratic Assignment Problem, Combinatorial Optimization, Graph Neural Networks.

✦

## 1 INTRODUCTION AND PRELIMINARIES

Graph matching (GM) is a fundamental problem which is NP-complete in general [1]. It has various applicability and connection with vision and learning, which involves establishing node correspondences between two graphs based on the node-to-node and edge-to-edge affinity [2], [3]. This differs from the point-based techniques e.g. RANSAC [4] and iterative closest point (ICP) [5] without considering edge information.

We start with two-graph matching, which can be written as quadratic assignment programming (QAP) [6], where $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ is a (partial) permutation matrix encoding node-to-node correspondence (with constraints in second line of Eq. (1)), and vec($\mathbf{X}$) is its column-vectorized version:

$$J(\mathbf{X}) = \text{vec}(\mathbf{X})^\top \mathbf{K} \text{vec}(\mathbf{X})$$
$$s.t. \quad \mathbf{X} \in \{0,1\}^{n_1 \times n_2}, \mathbf{X}\mathbf{1}_{n_2} = \mathbf{1}_{n_1}, \mathbf{X}^\top \mathbf{1}_{n_1} \leq \mathbf{1}_{n_2} \tag{1}$$

Here $\mathbf{1}_n$ means column vector of length $n$ whose elements all equal to 1, and $\mathbf{K} \in \mathbb{R}^{n_1 n_2 \times n_1 n_2}$ is the so-called affinity matrix [7]. Its diagonal and off-diagonal elements store the node-to-node and edge-to-edge affinities. For graph matching, the objective $J(\mathbf{X})$ is maximized, assuming perfect matching corresponds to the highest affinity score. One

- *R. Wang, J. Yan, and X. Yang are with Department of Computer Science and Engineering, and MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai, 200240, P.R. China. E-mail: {runzhong.wang, yanjunchi, xkyang}@sjtu.edu.cn Junchi Yan is the correspondence author.*

*Manuscript, under review.*

TABLE 1
Summary of existing learning-free and learning-based methods on two popular formulations of QAP. Note that Lawler's QAP (Eq. 1) is most general which incorporates the Koopmans-Beckmann's QAP (Eq. 3).

|  | Koopmans-Beckmann's QAP | Lawler's QAP |
|---|---|---|
| learning-free | [9], [10], [11], [12] | [2], [7], [13], [14], [15] |
| learning-based | [16], [17] | ours |

popular embodiment of $\mathbf{K}$ is fixed Gaussian kernel with Euclid distance over pairs of edge features [2], [8]:

$$K_{ia,jb} = \exp\left(\frac{||\mathbf{f}_{ij} - \mathbf{f}_{ab}||^2}{\sigma^2}\right) \tag{2}$$

where $\mathbf{f}_{ij}$ is the feature vector of the edge $E_{ij}$ in graph 1, $\mathbf{f}_{ab}$ from edge $E_{ab}$ in graph 2. $K_{ia,jb}$ is indexed by $(an_1 + i, bn_1 + j)$ under its matrix form.

Note Eq. (1) in literature is called Lawler's QAP [18], which can incorporate other special forms. For instance, the popular Koopmans-Beckmann's QAP [19] is written by:

$$J(\mathbf{X}) = \text{tr}(\mathbf{X}^\top \mathbf{F}_1 \mathbf{X} \mathbf{F}_2) + \text{tr}(\mathbf{K}_p^\top \mathbf{X}) \tag{3}$$

where $\mathbf{F}_1 \in \mathbb{R}^{n_1 \times n_1}$, $\mathbf{F}_2 \in \mathbb{R}^{n_2 \times n_2}$ are weighted adjacency matrices. $\mathbf{K}_p \in \mathbb{R}^{n_1 \times n_2}$ is node-to-node affinity matrix. Its connection to Lawler's QAP becomes clear by letting $\mathbf{K} = \mathbf{F}_2 \otimes_{\mathcal{K}} \mathbf{F}_1$ ($\otimes_{\mathcal{K}}$ means Kronecker product).

As shown in Tab. 1, traditional learning-free solvers have been extensively studied for both QAP formulations in Eq. (1, 3), and there exist some recent advances in learning
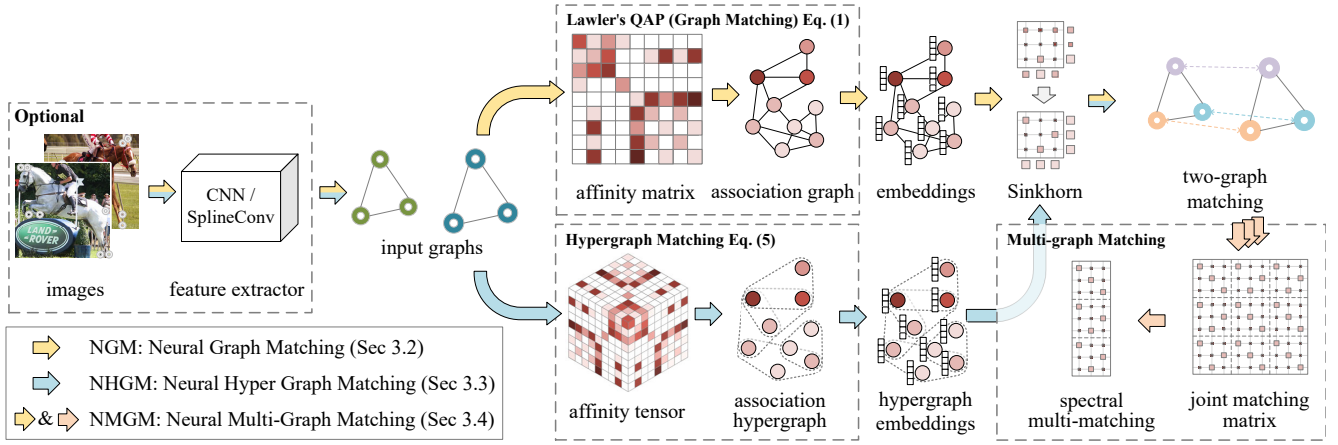
Fig. 1. Overview of the proposed neural graph matching pipeline. The method directly handles Lawler's QAP based on the embedding on the association graph. We further extend it to hypergraph matching by replacing the association graph with an association hypergraph (see Sec. 3.3), as well as to multiple graph matching by differentiable spectral multi-matching (see Sec. 3.4).

Koopmans-Beckmann's QAP [16], [17]. In this paper, we propose the first learning-based algorithm tackling the most general QAP form – Lawler's QAP, and show its generalization to higher-order and multi-graph scenarios.

The above QAP models involve the second-order affinity, and can also be generalized to the higher-order case. A line of works [20], [21], [22], [23] adopt tensor marginalization based model for $m$-order ($m \geq 3$) hypergraph matching, resulting in a higher-order assignment problem:

$$J(\mathbf{x}) = \mathbf{H} \otimes_1 \mathbf{x} \otimes_2 \mathbf{x} \ldots \otimes_m \mathbf{x}$$
$$s.t. \quad \mathbf{X}\mathbf{1}_{n_2} = \mathbf{1}_{n_1}, \mathbf{X}^\top \mathbf{1}_{n_1} \leq \mathbf{1}_{n_2} \tag{4}$$

where $\mathbf{x} = \text{vec}(\mathbf{X}) \in \{0,1\}^{n_1 n_2 \times 1}$ is the column-vectorized form, and $\mathbf{H}$ is the $m$-order affinity tensor whose $(n_1 n_2)^m$ elements record the affinity between two hyperedges, operated by tensor product $\otimes_k$ [24]:

$$(\mathbf{H} \otimes_k \mathbf{x})_{\ldots,i_{k-1},i_{k+1},\ldots} = \sum_{i_k=1}^{n_1 n_2} \mathbf{H}_{\ldots,i_{k-1},i_k,i_{k+1},\ldots} \cdot \mathbf{x}_{i_k} \tag{5}$$

where $\otimes_k$ can be regarded as tensor marginalization at dimension $k$. Details of tensor multiplication can be referred to Sec. 3.1 in [21]. Most existing hypergraph matching works assume the affinity tensor is invariant w.r.t. the index of the hyperedge pairs for computational tractability.

As discussed above, either graph matching or hypergraph matching problem involves solving a combinatorial optimization problem. However, the objective functions may be biased and even the mathematically optimal solution can depart from the perfect matching in reality, either due to the noise observation or limited modeling capacity, or both. In fact, traditional methods are mostly based on pre-defined shallow affinity with limited capacity e.g. Gaussian kernel with Euclid distance (see Eq. (2)), which has difficulty in providing enough flexibility for real-world data. This issue is partially addressed by affinity-learning based graph matching [17], [25], [26]. Along this promising direction, in this paper, a novel network based solver is proposed to directly learn Lawler's QAP whereby the affinity learning is also incorporated, as shown in Fig. 1. This approach is further extended to the case of joint matching of multiple

graphs, which has been an important scenario in practice and has received wide attention in literature [8], [27], [28], [29], [30] however learning has not been considered. Also hypergraph matching is enabled in our framework.

Specifically, the proposed matching nets consist of several learnable layers as detailed in Fig. 4: 1) CNN layers taking raw images for node (and edge) feature extraction; 2) Spline convolution layers encoding geometric features to node (and edge) features; 3) affinity metric learning layer for generating the affinity matrix i.e. the association graph; 4) vertex embedding layers using the association graph as input for vertex classification; 5) Sinkhorn net to convert the vertex score matrix into doubly-stochastic matrix. Sinkhorn technique is also adopted in the embedding module to introduce matching constraints; 6) cross-entropy loss layer whose input is the output of the Sinkhorn layer.

Note that the first two components are optional and can be treated as a plugin in the pipeline, and have nothing to do with Lawler's QAP. In contrast, the peer network [25] only allows for learning of node CNN features on images and their similarity metric i.e. the component 1) and 2). In fact, [25] is inapplicable to learning the QAP model. Our embedding also differs from [17], [31], [32] as raw individual graphs must be required for embedding in these works. In fact, [16] shows that embedding on individual graphs can deal with some special cases of Koopmans-Beckmann's QAP, which is also a special case of Lawler's QAP as discussed above. Direct learning on Lawler's QAP enables a learning-based solver for real-world combinatorial problems beyond vision, e.g. QAPLIB problem instances [33], which can not be readily handled by previous graph matching learning algorithms.

Furthermore, we devise two generalizations to the above matching network. 1) hypergraph matching by embedding a higher-order affinity tensor; 2) multiple graph matching by devising an end-to-end compatible matching synchronization module by using the popular spectral fusion technique [29], [34]. The performance can also be boosted by adopting edge-embedding layers.

The highlights of this paper are summarized as follows:
i) We show how to develop a deep network to directly

TABLE 2
Summary of existing literature in learning graph matching based on different types of assignment problems solved by learning, learning modules including CNN, GNN and affinity metric, where GNN embedding is performed and loss functions. KB-QAP abbreviates Koopmans-Beckmann's QAP in Eq. (3). Lawler's QAP in Eq. (1) is the most general QAP form and higher-order assignment in Eq. (4) is its higher-order extension.

| method | learned neural net solver | multiple-graph | CNN | GNN module | embedded graph | affinity metric | loss function |
|---|---|---|---|---|---|---|---|
| Nowak *et al.* [16] | special case of KB-QAP | none | none | GCN | individual graphs | inner-product | multi-class cross-entropy |
| GMN [25] | none | none | VGG16 | none | none | weighted exponential | pixel offset regression |
| Zhang *et al.* [31] | none | none | none | message-passing based CMPNN | individual graphs | inner-product | multi-class cross-entropy |
| PCA-GM [17] & IPCA-GM [35] | special case of KB-QAP | none | VGG16 | GCN+cross-graph convolution | individual graphs | weighted exponential | binary cross-entropy |
| CIE-H [36] | special case of KB-QAP | none | VGG16 | edge embedding+ cross-graph conv. | individual graphs | weighted exponential | binary cross-entropy with Hungarian attention |
| LCS [32] | special case of Lawler's QAP | none | VGG16 | GCN | association graph | fully-connected neural network | binary cross-entropy |
| BBGM [37] | none | unlearned fixed solver [38] | VGG16 | SplineConv | individual graphs | weighted inner-product | Hamming distance |
| NGM (ours) | **general case of Lawler's QAP (most general)** | none | VGG16 | matching-aware GCN | association graph | weighted exponential | binary cross-entropy |
| NHGM (ours) | **higher-order assignment** | none | VGG16 | matching-aware hyper-GCN | association hyper-graph | weighted exponential | binary cross-entropy |
| NMGM (ours) | **general case of Lawler's QAP (most general)** | **end-to-end spectral method** | VGG16 | matching-aware GCN | association graph | weighted exponential | binary cross-entropy |
| NGM-v2 (ours) | **general case of Lawler's QAP (most general)** | none | VGG16 | SplineConv+ match-aware GCN | individual graphs +asso. graph | weighted inner-product | binary cross-entropy |
| NHGM-v2 (ours) | **higher-order assignment** | none | VGG16 | SplineConv+match-aware hyper-GCN | individual graphs +asso. hyper-graph | weighted inner-product | binary cross-entropy |
| NMGM-v2 (ours) | **general case of Lawler's QAP (most general)** | **end-to-end spectral method** | VGG16 | SplineConv+ match-aware GCN | individual graphs +asso. graph | weighted inner-product | binary cross-entropy |

tackle the (most) general graph matching formulation i.e. Lawler's Quadratic Assignment Problem beyond vision problems, in the sense of allowing the affinity matrix as the raw input. This is fulfilled by regarding the affinity matrix as an association graph, whose vertices can be embedded by a deep graph neural network (GNN) [39] for classification, with a novel matching-aware graph convolution scheme. In contrast, existing works [17], [25], [26], [31] start with individual graphs' node and edge features for affinity learning instead of pairwise affinity encoded in affinity matrix.

ii) Our network solver for Lawler's QAP can be trained either in a supervised setting given ground truth node correspondence from labeled training set (e.g. for image matching), or by the final matching score without supervision (e.g. for QAPLIB problems).

iii) We extend our second-order graph matching networks to the hypergraph (third-order) matching case. This is fulfilled by building the hyperedge based association hypergraph to replace the second-order one. To our best knowledge, this is the first work for deep learning of hypergraph matching (with explicit treatment on the hyperedges).

iv) We also extend our matching network to the multiple-graph matching case by end-to-end spectral multi-graph matching, with explicit treatment for stabilized learning. To our best knowledge, there is no end-to-end multiple-graph matching neural network in the existing literature.

v) Experimental results on synthetic and real-world data show the effectiveness of our devised components. The extended versions for hypergraph matching and multiple-graph matching also show competitive performance. Our model can learn with the Lawler's QAP as input while state-of-the-art graph matching networks [17], [25], [31] cannot. This allows for the evaluation of our network on the QAPLIB benchmark directly, which to our best knowledge, is the first test for network-based methods for QAPLIB.

The paper goes as follows. Section 2 discusses the related work to graph matching and its recent learning based methods. Section 3 describes the main approach and in Section 4

we discuss the experiments. Section 5 concludes this paper. Source code is made public available on the project page.

## 2 RELATED WORK

We mainly discuss related works and techniques on learning graph matching [40]. Readers are referred to the survey [41] for an enlarged overview of the topic of graph matching. Also, a more broad perspective on learning for combinatorial optimization is referred to [42].

### 2.1 Learning-free Graph Matching Methods

**Two-graph matching and QAP.** Lawler's Quadratic Assignment Problem [18] is known for its application of matching two graphs by maximizing a quadratic objective function. Traveling salesman problem (TSP) and Koopmans-Beckmann's QAP are two popular variants from Lawler's QAP, with their wide range of application beyond vision, e.g. economic activities modeled by Koopmans-Beckmann's QAP [19], [33]. The most general Lawler's QAP also refers to two-graph matching in pattern recognition, and is traditionally addressed in a learning-free setting. Classically, small or medium sized problems are tractable by branch-and-bound with dual bound solved approximately [9], [13]. Modern approximate solvers [2], [7], [14], [15], [43] achieve better accuracy-speed trade-off and thus are more applicable to larger-sized problems. In two-graph matching, most methods focus on seeking approximate solution given fixed affinity model which is often set in simple parametric forms. Euclid distance in node/edge feature space together with a Gaussian kernel to derive a non-negative similarity, is widely used in the above works.

**Hypergraph matching methods.** Going beyond the traditional second-order graph matching, hypergraphs have been built for matching [23] and their affinity is usually represented by a tensor to encode the third-order [22], [24], [44] or even higher-order information [45]. The advantage is that the model can be more robust against noise at the cost of exponentially increased complexity for both time and space.
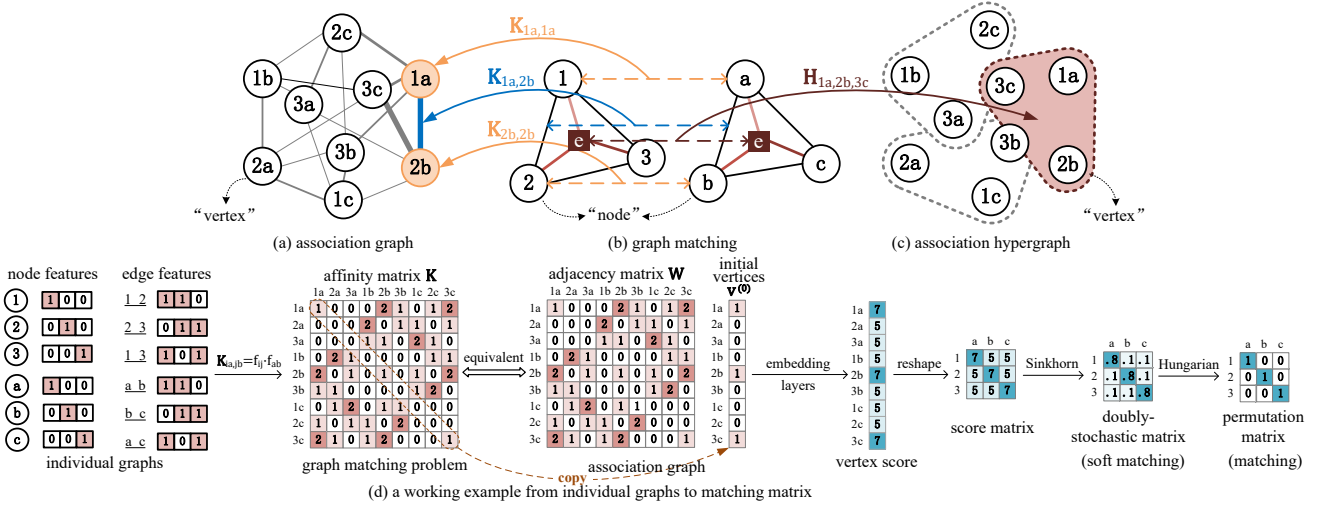
Fig. 2. Graphs with affinity matrix $\mathbf{K}$ and affinity tensor $\mathbf{H}$ in (b), w.r.t (a) association graph and (c) association hypergraph. We distinguish the nodes on association (hyper)graph and individual graphs for matching by terms *vertex* and *node* through this paper. The node-to-node matching problem in (b) can therefore be formulated as the vertex classification task on the association graph whose edge weights can be induced by the affinity matrix. Such a perspective is also widely taken in literature for graph matching e.g. [2] and hypergraph matching [24]. (d) shows a toy working example from individual graphs to final matching result: affinity matrix $\mathbf{K}$ is built from individual graphs, and the matching problem is equivalent to vertex classification on the association graph. Matching-aware embedding and vertex classification are applied on association graph to generate vertex scores, followed with reshaping and Sinkhorn normalization to obtain a double-stochastic matrix i.e. a convex hull of permutation matrix.

**Multiple-graph matching methods.** It has been recently actively studied for its practical utility against local noise and ambiguity. The hope is that the joint matching of multiple graphs can provide a better venue to fuse the information across graphs, leading to better robustness against local noise and ambiguity. Among the literature, a thread of works [28], [29] first generate the pairwise matching between two graphs via certain two-graph matching solvers, and then impose cycle-consistency on the pairwise matchings to improve the matching accuracy. The other line of methods impose cycle-consistency during the iterative finding of pairwise matchings and usually can achieve better results [8], [27], [46], [47]. The online setting for solving multiple graph matching is studied in [48].

Note both hypergraph or multiple graph matching paradigms try to improve the affinity model either by lifting the affinity order or imposing additional consistency regularization. As shown in the following, another possibly more effective and efficient way is adopting learning to find more adaptive affinity model parameters, or further improving the solver with deep neural networks.

## 2.2 Learning-based Graph Matching Methods

**Non-deep learning methods.** The structural SVM based supervised learning method [26] incorporates earlier graph matching learning methods [49], [50], [51]. Learning can also be fulfilled by unsupervised [51] and semi-supervised [52]. In these earlier works, no neural network is adopted until the recent seminal work [25].

**Deep-learning methods.** A pioneer work [16] considers the alignment of graphs by embedding on individual graphs, which can be regarded a special case of Koopmans-Beckmann's QAP. Deep learning is recently applied for graph matching on images [25], whereby convolutional neural network (CNN) is used to extract node features from images followed with spectral matching and CNN

is learned using a regression-like node correspondence supervision. This work is improved by introducing GNN to encode structural [17], [35] or geometric [31] information, with a combinatorial loss based on cross-entropy loss, and Sinkhorn network [53] as adopted in [17]. Yu *et al.* [36] extends [17] by edge embedding and Hungarian-based attention mechanism to stabilize end-to-end training.

We also note one recent important work on learning for graph matching, namely Learning Combinatorial Solver (LCS) [32]. It follows the line of research [2], [7], [54] in building an association graph from input images and solves graph matching by vertex classification on association graph. The discussion in [32] is basically restricted to the image matching setting. Though there is much space to generalize their work to our case, while no explicit scheme is given in terms of addressing the most general Lawler's QAP form. Meanwhile, it is unclear in their method for how to incorporate the matching constraint in the node scoring procedure, which has been addressed by our devised Skinhorn embedding. In fact, our method is directly motivated by developing a general Lawler's QAP solver, beyond image matching. It also enjoys the flexibility of readily adopting better feature extractors e.g. [37] or those beyond vision. Moreover, our work further develops hypergraph matching and multiple graph matching under the neural learning framework, which are new in literature. All these features are not well explored by LCS [32].

As summarized in Tab. 2 concerning deep neural network-based graph matching algorithms, one shortcoming of existing graph matching networks is that they cannot directly deal with the most general Lawler's QAP form which limits their applicability to tasks when no individual graph information is available (see QAPLIB – http://anjos.mgi.polymtl.ca/qaplib/). In contrast, our method can directly work with the affinity matrix, and we further extend to dealing with affinity tensor for hypergraph

matching, as well as the setting under multiple graphs.

## 3 PROPOSED APPROACHES

In Sec. 3.1, we introduce the preliminary concept association graph, on which our methods are based. In Sec. 3.2, we present Neural Graph Matching (NGM) network, which can solve Lawler's QAP for two-graph matching directly. Also, we show the extension to hypergraph matching, i.e. Neural Hyper-Graph Matching (NHGM) in Sec. 3.3, and to multiple graph matching i.e. Neural Multi-Graph Matching (NMGM) in Sec. 3.4. All these three settings to our knowledge have been hardly addressed by neural network solvers before. In Sec. 3.5 we introduce NGM/NHGM/NMGM-v2 with significantly improved image matching accuracy by exploiting enhanced feature extractor.

### 3.1 Preliminaries

Our models aim to match weighted graph $G^1 = (V^1, E^1)$ and $G^2 = (V^2, E^2)$ (in capital letters), where the superscript means the index of graphs and the subscript means the index of nodes. Without loss of generality, $|V^1| = n_1 = n_{in}$ are all inlier nodes, and $|V^2| = n_2 = n_{in} + n_{out}$ contains both inliers and optional outliers. $E^1, E^2$ are attributed edge sets with second-order features in graphs and $|E^1| = n_{e1}, |E^2| = n_{e2}$. Lawler's QAP as given in Eq. (1) is relaxed via popular doubly-stochastic relaxation:

$$J(\mathbf{S}) = \text{vec}(\mathbf{S})^\top \mathbf{K} \text{vec}(\mathbf{S}), \qquad (6)$$
$$\mathbf{S} \in [0,1]^{n_1 \times n_2}, \quad \mathbf{S}\mathbf{1}_{n_2} = \mathbf{1}_{n_1}, \quad \mathbf{S}^\top \mathbf{1}_{n_1} \leq \mathbf{1}_{n_2}$$

where $\mathbf{S}$ is a (partial) doubly-stochastic matrix, where all its rows sum to 1 and all its column sums are $\leq 1$. For affinity matrix $\mathbf{K} \in \mathbb{R}^{n_1 n_2 \times n_1 n_2}$, diagonal elements $\mathbf{K}_{ia,ia} = \mathbf{s}_v(V_i^1, V_a^2)$ are first order (node) similarities and off-diagonal elements $\mathbf{K}_{ia,jb} = \mathbf{s}_e(E_{ij}^1, E_{ab}^2)$ are second order (edge) similarities, where $\mathbf{s}_v, \mathbf{s}_e$ are similarity measurements for nodes and edges, respectively.

As shown in Fig. 2, graph matching can be viewed in a perspective based on the definition of the so-called association graph $\mathcal{G}^A = (\mathcal{V}^A, \mathcal{E}^A)$ [2], [7] (in handwritten letters with superscript $A$). **To avoid ambiguity between graphs and association graphs, we name the entities in graphs ($V$) as *nodes* and the entities in association graph ($\mathcal{V}^A$) as *vertices*.** Readers should distinguish these two concepts as they will be repeatedly encountered through this paper.

The vertices of association graph $\mathcal{V}^A = V^1 \times V^2$ encode candidate node-to-node correspondence $\mathcal{V}_{ia}^A = (V_i^1, V_a^2)$ corresponding to the matching matrix $\mathbf{X}_{i,a}$, therefore the vectorized assignment matrix $\text{vec}(\mathbf{X})$ is equivalent to the vertex set of association graph. The edges $\mathcal{E}^A$ represent the agreement between two pairs of correspondence $\mathcal{E}_{ia,jb}^A = \{(V_i^1, V_a^2), (V_j^1, V_b^2)\}$ modeled by $\mathbf{K}_{ia,jb}$, so that the off-diagonal part of affinity matrix $\mathbf{K}$ is equivalent to the adjacency matrix of association graph. The matching between two graphs can therefore be transformed into vertex classification on the association graph, following [2], [7]. In this paper, diagonal elements $\mathbf{K}_{ia,ia}$ are further assigned as vertex attributes $\mathcal{V}^A$, to better exploit the first-order similarities. Such formulation can also be generalized to hypergraph matching problems, with edges replaced by hyperedges, as



(a) graph

(b) adjacency matrix $= \bar{\mathbf{G}} \, \bar{\mathbf{H}}^\top$

(c) $\bar{\mathbf{G}}$

(d) $\bar{\mathbf{H}}$

(e) incidence matrix $= \bar{\mathbf{H}} \cdot \bar{\mathbf{G}}$
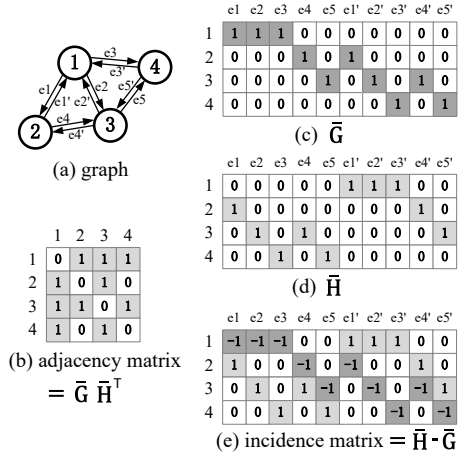
Fig. 3. A toy example of connectivity matrices (c) $\bar{\mathbf{G}}$, (d) $\bar{\mathbf{H}}$, and their connection to (a) the original graph, (b) adjacency matrix and (e) incidence matrix. All edges are directed, and $\bar{\mathbf{G}}_{i,k} = \bar{\mathbf{H}}_{j,k} = 1$ means edge $k$ starts from node $i$ and ends at node $j$. In incident matrix, -1 denotes the starting node and 1 denotes the ending node of all edges.

shown in Fig. 2(c). The association graph prohibits links that violate the one-to-one matching constraint (e.g. there is no link between vertex '1a' and '1c' in Fig. 2(a)).

### 3.2 NGM: Neural Graph Matching for QAP

Based on the formulations and association graph introduced in Sec. 3.1, our proposed Neural Graph Matching (NGM) solves relaxed Lawler's QAP in Eq. (6), by vertex classification via Graph Convolutional Networks (GCN) [55] with novel matching-aware embedding modules. The vertex classification is performed on the association graph induced by the affinity matrix, followed by a Sinkhorn operator. As shown by existing learning-free graph matching solvers [2], [7], graph matching problem is equivalent to vertex classification on the association graph. NGM accepts either raw image (with jointly learned CNN and affinity metric), or affinity matrix (without CNN or affinity metric), and learns end-to-end from ground truth correspondence or by pure optimization for QAPLIB problems.

#### 3.2.1 Affinity matrix building from natural images

Our graph matching net allows either affinity matrix or raw images as input. The image processing module is optional and treated as a plug-in for dealing with images, following the protocol in [25] whereby the affinity matrix is built from pre-given keypoints in images. As shown in the upper half of Fig. 4, image features are extracted by learnable CNN layers such as VGG16 [56]. Given two input images with labeled keypoints, we adopt CNN layers to extract per-node features $\bar{\mathbf{F}}^1, \bar{\mathbf{U}}^1 \in \mathbb{R}^{n_1 \times d}$ for $G^1$ and $\bar{\mathbf{F}}^2, \bar{\mathbf{U}}^2 \in \mathbb{R}^{n_2 \times d}$ for $G^2$, where $d$ is feature dimension size and $\bar{\mathbf{F}}, \bar{\mathbf{U}}$ are extracted from different CNN layers (e.g. VGG16 `relu5_1` for $\bar{\mathbf{F}}$ and `relu4_2` for $\bar{\mathbf{U}}$) and utilized for edge representation and node representation, respectively. Features are obtained by bi-linear interpolation on the CNN feature map. As shown in Fig. 3, the connectivity of two graphs are represented by $\bar{\mathbf{G}}^1, \bar{\mathbf{H}}^1 \in \{0,1\}^{n_1 \times n_{e1}}$ and $\bar{\mathbf{G}}^2, \bar{\mathbf{H}}^2 \in \{0,1\}^{n_2 \times n_{e2}}$, where $\bar{\mathbf{A}}^1 = \bar{\mathbf{G}}^1 \bar{\mathbf{H}}^{1\top}, \bar{\mathbf{A}}^2 = \bar{\mathbf{G}}^2 \bar{\mathbf{H}}^{2\top}$ are the adjacency matrices
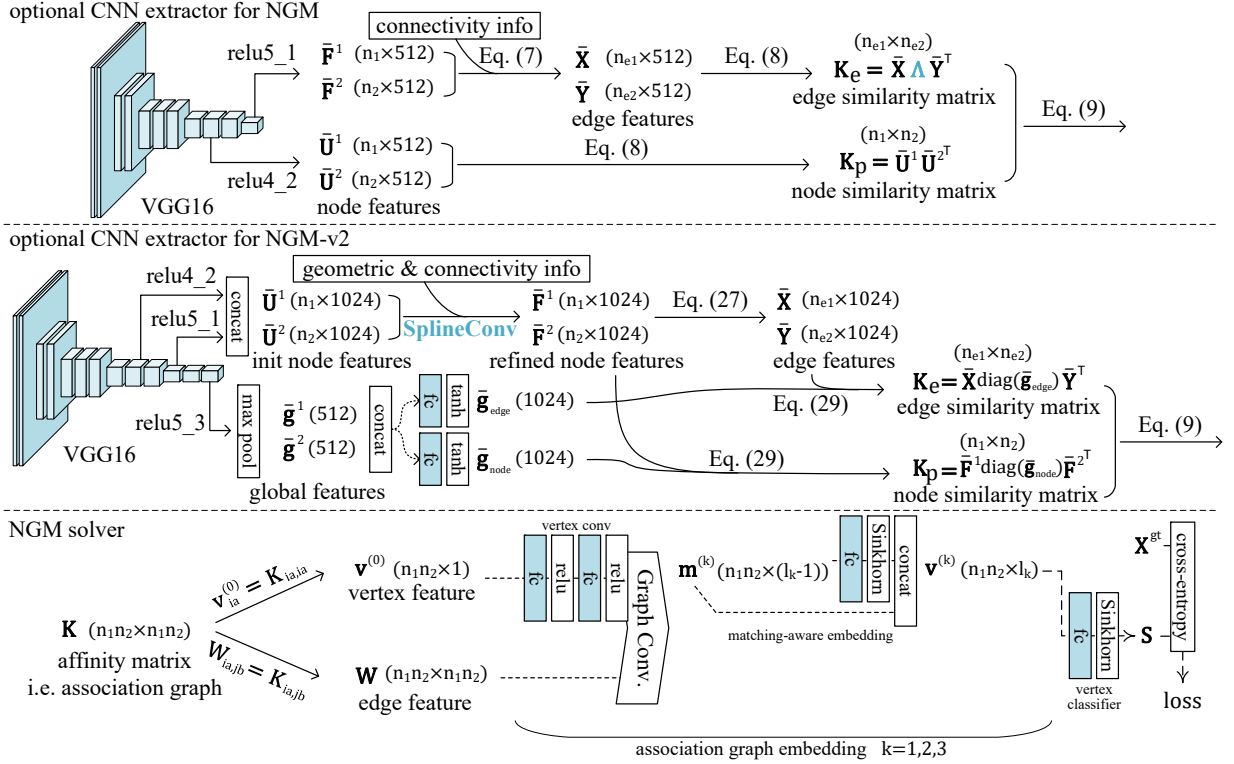
Fig. 4. The proposed NGM & NGM-v2 architecture for two-graph matching. The components with blue FC layers i.e. vertex convolution, matching-aware embedding module and vertex classifier are jointly learned with optional CNN, SplineConv (in NGM-v2) and similarity metric $\mathbf{\Lambda}$ (in NGM).

of two graphs, and $\bar{\mathbf{G}}_{i,k} = \bar{\mathbf{H}}_{j,k} = 1$ means edge $k$ links node $i$ to node $j$. The edge representations are built by concatenating node features at both ends of the edge:

$$\bar{\mathbf{X}} = [\bar{\mathbf{G}}_1^\top \bar{\mathbf{F}}_1 \quad \bar{\mathbf{H}}_1^\top \bar{\mathbf{F}}_1], \ \bar{\mathbf{Y}} = [\bar{\mathbf{G}}_2^\top \bar{\mathbf{F}}_2 \quad \bar{\mathbf{H}}_2^\top \bar{\mathbf{F}}_2] \quad (7)$$

where $[\ \cdot \ \cdot\ ]$ means concatenating two matrices along columns. The node-to-node similarity matrix $\mathbf{K}_p \in \mathbb{R}^{n_1 \times n_2}$ and edge-to-edge similarity $\mathbf{K}_e \in \mathbb{R}^{n_{e1} \times n_{e2}}$ are built via

$$\mathbf{K}_e = \bar{\mathbf{X}} \mathbf{\Lambda} \bar{\mathbf{Y}}^\top, \mathbf{K}_p = \bar{\mathbf{U}}^1 \bar{\mathbf{U}}^{2\top} \quad (8)$$

where $\mathbf{\Lambda} \in \mathbb{R}^{2d \times 2d}$ is the learnable parameter for affinity metric. The QAP affinity matrix is built following the factorized formulation of $\mathbf{K}$ [57]:

$$\mathbf{K} = \mathrm{diag}(\mathrm{vec}(\mathbf{K}_p)) + (\bar{\mathbf{G}}_2 \otimes_\mathcal{K} \bar{\mathbf{G}}_1)\mathrm{diag}(\mathrm{vec}(\mathbf{K}_e))(\bar{\mathbf{H}}_2 \otimes_\mathcal{K} \bar{\mathbf{H}}_1)^\top \quad (9)$$

where $\mathrm{diag}(\cdot)$ means building a diagonal matrix from input vector, and $\otimes_\mathcal{K}$ means Kronecker product. All the forementioned operations allow back propagation, and we adopt the efficient GPU implementation provided by [17].

### 3.2.2 Association graph construction

We derive the association graph from the affinity matrix $\mathbf{K}$. The weighted adjacency matrix of association graph $\mathbf{W}$ comes from the off-diagonal elements of $\mathbf{K}$. We denote $\mathbf{v}^{(k)} \in \mathbb{R}^{n_1 n_2 \times l_k}$ as $l_k$-dimensional vertex embeddings on layer $k$ (starting with $k = 0$). The initial embeddings at $k = 0$ are scalar, i.e. $l_0 = 1$, taken from the diagonal of $\mathbf{K}$.

$$\mathbf{W}_{ia,jb} = \mathbf{K}_{ia,jb}, \quad \mathbf{v}_{ia}^{(0)} = \mathbf{K}_{ia,ia} \quad (10)$$

$\mathbf{W}$ contains both connectivity and weight information in the association graph. In case when the first-order similarity $\mathbf{K}_{ia,ia}$ is absent, we can assign a constant (e.g. 1) for all $\mathbf{v}^{(0)}$.

### 3.2.3 Matching aware embedding of association graph

The matching problem can be transformed to selecting the vertices in the association graph that encode the node-to-node correspondence between two input graphs, as shown in Fig. 2. For vertex classification on the association graph, we use GCN [55] for its effectiveness and simplicity.

We define the (unweighted) adjacency matrix of association graph $\mathbf{A} \in \{0,1\}^{n_1 n_2 \times n_1 n_2}$: $\mathbf{A}_{ia,jb} = 1$ if $\mathbf{K}_{ia,jb} > 0$ and otherwise 0. $\mathbf{A}_{ia,jb}$ serves as an indicator whether there exists an edge between vertices $ia$ and $jb$ in the association graph. Recall that the vertex $ia$ represents the matching between node $i$ and $a$ from separate input graphs (see Fig. 2). In association graph, an edge exists between $ia$ and $jb$ if and only if the node-to-node matchings $i$ to $a$ and $j$ to $b$ can co-exist, and there exists an edge-to-edge affinity score defined between edges $ij$ and $ab$. Since $\mathbf{A}$ is symmetric, we compute the degree matrix for normalization:

$$\mathbf{D} = \mathrm{diag}(\mathbf{A} \mathbf{1}_{n_1 n_2}) \quad (11)$$

where $\mathrm{diag}(\cdot)$ builds a diagonal matrix from input vector. The vertex aggregation step is according to:

$$\mathbf{m}^{(k)} = \mathbf{D}^{-1} \mathbf{W} f_m(\mathbf{v}^{(k-1)}) + f_v(\mathbf{v}^{k-1}), \ \mathbf{v}^{(k)} = \mathbf{m}^{(k)} \quad (12)$$

where the message passing function $f_m : \mathbb{R}^{l_{k-1}} \to \mathbb{R}^{l_k}$ and vertex's self update function $f_v : \mathbb{R}^{l_{k-1}} \to \mathbb{R}^{l_k}$ are both implemented by networks with two fully-connected layers and ReLU activation.

The above general vanilla vertex embedding procedure in Eq. (12) does not consider the one-to-one assignment constraint for matching. Here we develop a matching constraint aware embedding model: in each layer a soft permutation (i.e. doubly-stochastic matrix) is scored via classifier with Sinkhorn network Classifier : $\mathbb{R}^{n_1 n_2 \times l_k} \to [0,1]^{n_1 \times n_2}$ (see discussions in Sec. 3.2.4) followed by vectorization operator $\text{vec}(\cdot)$. The predicted soft permutation is concatenated to vertex embeddings whereby matching information is considered in embedding layers. With $f_m, f_v : \mathbb{R}^{l_{k-1}} \to \mathbb{R}^{(l_k - 1)}$, such a matching-aware embedding scheme is denoted as **Sinkhorn embedding** in the rest of the paper.

$$\mathbf{v}^{(k)} = [\mathbf{m}^{(k)} \quad \text{vec}(\text{Classifier}(\mathbf{m}^{(k)}))] \quad (13)$$

where $[\;\cdot\quad\cdot\;]$ means concatenation. We experiment both vanilla vertex embedding in Eq. (12) and matching-aware Sinkhorn embedding in Eq. (13), to validate the necessity of adding assignment constraint.

### 3.2.4 Vertex classification with Sinkhorn network

As graph matching is equivalent to vertex classification on association graph (see Fig. 2), a vertex classifier with Sinkhorn network is adopted to predict the matching result. Specifically we use a single layer fully-connected classifier denoted by $f_c : \mathbb{R}^{l_k} \to \mathbb{R}$, followed by exponential activation with regularization factor $\alpha$:

$$\mathbf{s}_{ia}^{(k)} = \exp\left(\alpha f_c(\mathbf{v}_{ia}^{(k)})\right) \quad (14)$$

After reshaping classification scores into $\mathbb{R}^{n_1 \times n_2}$, one-to-one assignment constraint is enforced to $\mathbf{s}$ by Sinkhorn network [53], [58]. It takes a non-negative square matrix as input and outputs a doubly-stochastic matrix [53], [59]. As the scoring matrix can be non-square for different sizes of graphs, the input matrix $\mathbf{S} \in \mathbb{R}^{n_1 \times n_2}$ is padded into a square one (assume $n_1 \leq n_2$) with small elements e.g. $\epsilon = 10^{-3}$. A doubly-stochastic matrix is obtained by repeatedly running:

$$\mathbf{S} = \mathbf{S} \oslash (\mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top \cdot \mathbf{S}), \; \mathbf{S} = \mathbf{S} \oslash (\mathbf{S} \cdot \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top) \quad (15)$$

where $\oslash$ means element-wise division. By taking column-normalization and row-normalization in Eq. (15) alternatively, $\mathbf{S}$ converges to a doubly-stochastic matrix whose rows and columns all sum to 1. The dummy elements are discarded in the final output, whose column sum may be $< 1$ given umatched nodes from the bigger graph. Sinkhorn operator is fully differentiable and can be efficiently implemented by automatic differentiation techniques [60]. The proposed vertex classifier with Sinkhorn network is denoted as Classifier : $\mathbb{R}^{n_1 n_2 \times l_k} \to [0,1]^{n_1 \times n_2}$, which is also mentioned in matching-aware Sinkhorn embedding in Eq. (13).

### 3.2.5 Loss for end-to-end training

Recall the obtained predicted matrix $\mathbf{S}$ from the above procedure is a doubly-stochastic matrix. Each element can be regarded as a binary classification where each vertex should be classified to 1 (matched) or 0 (unmatched). Hence we adopt the binary cross-entropy as the final loss, given the ground truth node-to-node correspondence $\mathbf{X}^{gt}$:

$$\ell = -\sum_{i=1}^{n_1} \sum_{a=1}^{n_2} \mathbf{X}_{i,a}^{gt} \log \mathbf{S}_{i,a} + (1 - \mathbf{X}_{i,a}^{gt}) \log(1 - \mathbf{S}_{i,a}) \quad (16)$$
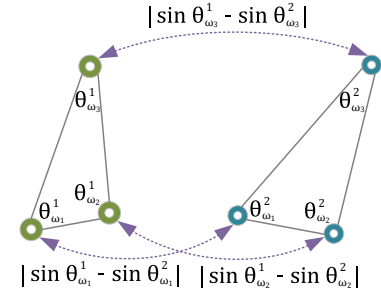


Fig. 5. Third-order affinity adopted in our hypergraph matching. It in general follows the previous hypergraph matching works [22], [24], by considering the similarity between two triangle's three inner angles.

Our approach also allows direct optimization over the affinity score objectives for QAPLIB problems, which will be discussed in Sec. 4.2 in details.

All the components are differentiable. Therefore, both NGM solver and the optional CNN feature extractor can be learned via backpropagation and gradient descent.

## 3.3 NHGM: Neural Hypergraph Matching

For Neural Hyper-Graph Matching (NHGM), higher-order structure is exploited for more robust matching. NHGM owns a nearly identical pipeline compared to NGM, while a more general message-passing scheme is devised for vertex classification in hypergraphs, as previously shown in [61], [62]. Due to the explosive computational cost ($O((n_1 n_2)^t)$ with order $t$), here we limit hypergraph to third-order which is also in line with the majority of existing works [22], [24], while the scheme is generalizable to any order $t$.

### 3.3.1 Association hypergraph construction

The second-order affinity matrix $\mathbf{K}$ is generalized to affinity tensor $\mathbf{H}^{\langle t \rangle}$ of order $t$ in hypergraph matching. In line with the hypergraph matching literature [21], [23], [24], the third-order affinity tensor is specified as:

$$\mathbf{H}_{\omega_1, \omega_2, \omega_3}^{\langle 3 \rangle} = \exp\left(-\left(\sum_{q=1}^{3} |\sin \theta_{\omega_q}^1 - \sin \theta_{\omega_q}^2|\right) / \sigma_3\right) \quad (17)$$

where $\theta_{\omega_q}^1, \theta_{\omega_q}^2$ denotes the angle in graph $\mathcal{G}_1$ and $\mathcal{G}_2$ of each correspondence $\omega_q$, respectively. An illustration of third order affinity can be found in Fig. 5, where the similarity between triplets of nodes is compared. Third order affinity is usually defined on geometric consistency and it preserves both scaling and rotation invariance.

Extending from the second-order association graph, a hyper association graph is built from $\mathbf{H}$. The association hypergraph $\mathcal{H}^A = (\mathcal{V}^A, \mathcal{E}^A)$ takes node-to-node correspondence $\omega = (V_i^1, V_j^2)$ as vertices $\mathcal{V}^A$ (which is consistent with second-order association graph) and higher-order similarity among $\{(V_{\omega_1}^1, V_{\omega_1}^2), \cdots, (V_{\omega_t}^1, V_{\omega_t}^2)\}$ as hyperedges $\mathcal{E}^A$, as shown by Fig. 2(c). Elements of $\mathbf{H}$ are adjacency weights for the association hypergraph accordingly. In NHGM, hypergraph convolution is defined for vertex classification.

### 3.3.2 Matching aware association hypergraph embedding

As an extension of Eq. (12), vertex embeddings are updated from all vertices linked by hyperedges in association hypergraph. We compute normalized degree tensor of order $t$:

$$\mathbf{D}_{\omega_1, \cdots, \omega_t}^{\langle t \rangle -1} = \mathbf{A}_{\omega_1, \cdots, \omega_t}^{\langle t \rangle} / (\mathbf{A}^{\langle t \rangle} \otimes_2 \mathbf{1} \cdots \otimes_t \mathbf{1})_{\omega_1} \quad (18)$$

Then an aggregation scheme extended from Eq. (12) is taken:

$$\begin{aligned} \mathbf{p}^{(k)} &= f_m^{\langle t \rangle}(\mathbf{v}^{(k-1)}), \ \mathbf{H}^{\langle t \rangle \prime} = (\mathbf{D}^{\langle t \rangle -1} \odot \mathbf{H}^{\langle t \rangle})_{t+1} \\ \mathbf{m}^{(k)} &= \sum_t \lambda_t \mathbf{H}^{\langle t \rangle \prime} \otimes_t \mathbf{p}^{(k)} \cdots \otimes_2 \mathbf{p}^{(k)} + f_v \\ \mathbf{v}^{(k)} &= \begin{bmatrix} \mathbf{m}^{(k)} & \text{vec}(\text{Classifier}(\mathbf{m}^{(k)})) \end{bmatrix} \end{aligned} \quad (19)$$

where $f_v$ abbreviates $f_v(\mathbf{v}^{(k-1)})$, $\otimes_i$ denotes tensor product by dimension $i$ (see Eq. (5)), $\odot$ denotes element-wise multiply, $(\cdot)_{t+1}$ means expanding along dimension $(t+1)$. $f_m^{\langle t \rangle} : \mathbb{R}^{l_{k-1}} \rightarrow \mathbb{R}^{(l_k - 1)}$ is message passing function at order $t$. Different orders of features are fused by weighted summation with $\lambda_t$.

The other modules of NHGM, including classifier and cross-entropy loss, are identical to NGM. Therefore, NGM can be viewed as a special case of NHGM, where the order is restricted to 2. The sparsity of $\mathbf{H}$ is exploited for both time and memory efficiency.

## 3.4 NMGM: Neural Multi-graph Matching

We explore learning multi-graph matching, where the information is fused among graphs by the so-called cycle-consistency. Cycle-consistency denotes a condition where the matching between any two graphs is consistent when passed through any other graphs, i.e. $\mathbf{X}_{ij} = \mathbf{X}_{ik}\mathbf{X}_{kj}$ for all $i, j, k$, which can be viewed as each graph matched to a $n$-sized reference $\mathbf{X}_i \in \{0, 1\}^{n \times n}$. The pariwise matching between $G_i, G_j$ can be represented with $\mathbf{X}_{ij} = \mathbf{X}_i \mathbf{X}_j^\top$. In this paper, we refer to the line of works [28], [29], [34], [63] involving post-synchronization given initial pairwise matchings. Spectral fusion is used in NMGM for its effectiveness and simplicity, and most importantly, its ability for end-to-end training. We assume all graphs are of equal size.

### 3.4.1 Building joint matching matrix

We first obtain initial two-graph matchings by NGM to build a symmetric joint matching matrix $\mathcal{S} \in \mathcal{R}^{nm \times nm}$. For each pair $G_i$ and $G_j$ with $n$ nodes, $\mathbf{S}_{ij} \in [0, 1]^{n \times n}$ is computed by NGM as the soft (i.e. doubly-stochastic) matching matrix. For $m$ graphs, $\mathcal{S}$ can be built from all combinations of pairwise matchings $\mathbf{S}_{ij}$:

$$\mathcal{S} = \begin{pmatrix} \mathbf{S}_{00} & \cdots & \mathbf{S}_{0m} \\ \vdots & \ddots & \vdots \\ \mathbf{S}_{m0} & \cdots & \mathbf{S}_{mm} \end{pmatrix} \quad (20)$$

where $\mathcal{S}$ is of size $mn \times mn$. For the diagonal part of $\mathcal{S}$, $\mathbf{S}_{ii}$ are all identical matrices. Note $\mathbf{S}_{ij}$ are all square matrices. The objective of spectral fusion results in getting a cycle-consistent joint matching matrix $\hat{\mathcal{S}}$, whose innerproduct with $\mathcal{S}$ is maximized:

$$\max_{\hat{\mathcal{S}}} = \text{tr}(\hat{\mathcal{S}}^\top \mathcal{S}) \quad (21)$$

To ensure the cycle-consistency, it holds $\hat{\mathbf{S}}_{ij} = \hat{\mathbf{S}}_{ik}\hat{\mathbf{S}}_{kj}$ for all elements in $\hat{\mathcal{S}}$. We may select an arbitrary $k$ as the reference, omitting $k$ in the subscript, $\hat{\mathcal{S}}$ can be decomposed as matching to the reference:

$$\hat{\mathbf{U}}\hat{\mathbf{U}}^\top = \hat{\mathcal{S}}$$

$$\text{where} \quad \hat{\mathbf{U}} = \begin{pmatrix} \hat{\mathbf{S}}_0 \\ \vdots \\ \hat{\mathbf{S}}_m \end{pmatrix} \quad (22)$$

under ideal condition where $\hat{\mathbf{S}}_i$ are all permutation matrices, each column of $\hat{\mathbf{U}}$ is linearly independent and $\hat{\mathbf{U}}/\sqrt{m}$ are $n$ eigenvectors of $\hat{\mathcal{S}}$ with eigenvalue $m$. The permutation constraint in $\hat{\mathbf{S}}_i$ is relaxed for computational feasibility and Eq. (21) results in a generalized Rayleigh problem and solved via spectral fusion, as shown follows.

### 3.4.2 Differentiable spectral fusion of pairwise matchings

Multi-graph matching information can be fused by eigenvector decomposition (i.e. spectral method) on $\mathcal{S}$. Based on generalized Rayleigh problem, given $n$ nodes in each graph, we extract the eigenvectors corresponding to the top-$n$ eigenvalues of symmetric matrix $\mathcal{S}$:

$$\mathbf{U}\,\mathbf{\Sigma}\,\mathbf{U}^\top = \mathcal{S} \quad (23)$$

where diagonal matrix $\mathbf{\Sigma} \in \mathbb{R}^{n \times n}$ contains top-$n$ eigenvalues and $\mathbf{U} \in \mathbb{R}^{mn \times n}$ are the $n$ corresponding eigenvectors. It has been shown that the computation of eigenvalues and eigenvectors are differentiable [64] which makes them fixed components in our end-to-end learning network pipeline. The fusion of the input $\mathcal{S}$ can be written as follows, which can be seen as a smoothed version maximizing $\text{tr}(\hat{\mathcal{S}}^\top \mathcal{S})$:

$$\hat{\mathcal{S}} = m\,\mathbf{U}\,\mathbf{U}^\top \quad (24)$$

The gradient of eigen decomposition in Eq. (23) is [64]:

$$\frac{\partial L}{\partial \mathcal{S}} = \mathbf{U}\left(4m\left(\mathbf{Y}^\top \odot \left(\mathbf{U}^\top \left(\frac{\partial L}{\partial \hat{\mathcal{S}}}\right)_{sym} \mathbf{U}\right)_{sym}\right)\right)\mathbf{U}^\top$$

$$\text{where} \quad \mathbf{Y}_{ij} = \begin{cases} 1/(\sigma_i - \sigma_j) & i \neq j \\ 0 & i = j \end{cases} \quad (25)$$

where $L$ denotes the loss, $\odot$ means element-wise multiplication, $\mathbf{A}_{sym} = (\mathbf{A} + \mathbf{A}^\top)/2$ and $\sigma_i = \mathbf{\Sigma}_{ii}$ is the $i$-th eigenvalue. According to this backward formulation, if there exist non-distinctive eigenvalues, i.e. $\sigma_i = \sigma_j$ for $i \neq j$, a numerical divided-by-zero error will be caused. This issue usually happens when cycle-consistency (i.e. $\mathbf{S}_{ij} = \mathbf{S}_{ik}\mathbf{S}_{kj}$) is already met in the input $\mathbf{S}$, under such circumstances the fused matching results are nearly identical to the original matchings. To avoid numerical issues, we assign $\hat{\mathbf{S}}_{ij} = \mathbf{S}_{ij}$ to bypass eigendecomposition if the minimum residual among top-$n$ eigenvalues is smaller than tolerance $\delta$, e.g. $10^{-4}$. This strategy is found effective to stabilize learning.

The final matching results are obtained by differentiable Sinkhorn network:

$$\bar{\mathbf{S}}_{ij} = \text{Sinkhorn}(\exp(\hat{\alpha}\,\hat{\mathbf{S}}_{ij})) \quad (26)$$

where the fused two-graph matching $\hat{\mathbf{S}}_{ij}$ is from $\hat{\mathcal{S}}$ and $\exp(\hat{\alpha} \cdot )$ performs regularization for Sinkhron. Cross-entropy loss in Eq. (16) is applied to each $\bar{\mathbf{S}}_{ij}$ for supervised learning, which is similar to the supervised two-graph matching case in the paper.

### 3.5 Improved Matching by Enhanced Feature Extractor

Since our proposed method deals with the most general Lawler's QAP, the baseline feature extractor in Sec. 3.2.1 can be replaced by other enhanced feature extraction techniques. In this section, we refer to the novel feature extractor proposed by [37] in replacement of the feature extractor in Sec. 3.2.1, introducing the family of enhanced models for image matching problems – NGM/NHGM/NMGM-v2. In Sec. 3.5.1 we discuss how to build enhanced affinity matrix, and in Sec. 3.5.2 we propose a way of building hypergraph affinity tensor for NHGM-v2.

#### 3.5.1 Enhanced graph matching feature

The authors of [37] propose an enhanced deep learning feature extractor for graph matching problem on images based on SplineConv [65] and weighted inner product, which can seamlessly fit into our pipeline. This enhanced feature extractor is summarized in the second row of Fig. 4.

With some abuse of the notations from Sec. 3.2.1, we introduce the feature extractor as follows. Given two input images, we extract node features from relu4_2 and relu5_1 of VGG16, and then concatenate them to form the feature matrices $\bar{\mathbf{U}}^1 \in \mathbb{R}^{n_1 \times d}, \bar{\mathbf{U}}^2 \in \mathbb{R}^{n_2 \times d}$, where $d$ is the feature dimension of the concatenated features. Then we adopt two SplineConv[1] [65] layers on $\bar{\mathbf{U}}^1, \bar{\mathbf{U}}^2$ to produce refined features $\bar{\mathbf{F}}^1 \in \mathbb{R}^{n_1 \times d}, \bar{\mathbf{F}}^2 \in \mathbb{R}^{n_2 \times d}$. SplineConv [65] is a powerful graph convolution operator exploiting B-Spline kernel, encoding geometric features into node features. Therefore, SplineConv is suitable for feature refinement on image matching datasets. Readers are referred to the original paper for details about SplineConv. The edge features are constructed as the difference of its two nodes in the feature space:

$$\bar{\mathbf{X}} = \bar{\mathbf{G}}_1^\top \bar{\mathbf{F}}_1 - \bar{\mathbf{H}}_1^\top \bar{\mathbf{F}}_1, \ \bar{\mathbf{Y}} = \bar{\mathbf{G}}_2^\top \bar{\mathbf{F}}_2 - \bar{\mathbf{H}}_2^\top \bar{\mathbf{F}}_2 \qquad (27)$$

where $\bar{\mathbf{G}}, \bar{\mathbf{H}}$ are the connectivity matrices in Fig. 3.

As shown in the second row of Fig. 4, the enhanced feature extractor also contains a branch producing global features by max-pooling over the output of relu5_3 layer. The pooled global features from two graphs are then concatenated, and passed to a fc layer followed with tanh activation, producing global features $\bar{\mathbf{g}} \in \mathbb{R}^d$:

$$\bar{\mathbf{g}} = \tanh(\mathrm{fc}([\bar{\mathbf{g}}^1 \quad \bar{\mathbf{g}}^2])) \qquad (28)$$

where $[ \cdot \ \cdot ]$ means concatenation, and we compute separate $\bar{\mathbf{g}}_{\text{node}}, \bar{\mathbf{g}}_{\text{edge}}$ for node and edge features, respectively.

Node and edge similarity matrices are constructed based on weighted inner-product whereby $\bar{\mathbf{g}}$ as the weight:

$$\mathbf{K}_e = \bar{\mathbf{X}} \operatorname{diag}(\bar{\mathbf{g}}_{\text{node}}) \bar{\mathbf{Y}}^\top, \mathbf{K}_p = \bar{\mathbf{F}}^1 \operatorname{diag}(\bar{\mathbf{g}}_{\text{edge}}) \bar{\mathbf{F}}^{2\top} \qquad (29)$$

---

1. SplineConv is called SplineCNN in its original paper [65]. To avoid the ambiguity of the term "CNN" which usually represents convolution on images while SplineCNN is convolution on graphs, we name it as SplineConv in this paper.

where $\operatorname{diag}(\bar{\mathbf{g}})$ means building a diagonal matrix by placing $\bar{\mathbf{g}}$ on it. After obtaining $\mathbf{K}_e, \mathbf{K}_p$, the affinity matrix can be built by Eq. (9), and the resulting QAP is readily solved by NGM/NMGM as discussed in Sec. 3.2 and Sec. 3.4.

#### 3.5.2 Enhanced hypergraph affinity for NHGM-v2

Based on the the features defined in Sec. 3.5.1, we further design the hyperedge affinity for NHGM-v2. The hypergraph affinity is inspired by Fig. 5, however in NHGM the hypergraph affinity is based on geometric features, and in NHGM-v2 we define hypergraph affinity on the high-dimensional feature space. We regard the plane in Fig. 5 as the feature space formed by $\bar{\mathbf{F}}^1$ and $\bar{\mathbf{F}}^2$, and each point in Fig. 5 represents a node with features. Following Eq. (17), the third order affinity tensor is defined as the differences of angles in the feature space:

$$\mathbf{H}^{\langle 3 \rangle}_{\omega_1, \omega_2, \omega_3} = \exp\left( - \left( \sum_{q=1}^{3} |\cos\theta^1_{\omega_q} - \cos\theta^2_{\omega_q}| \right) / \sigma_3 \right) \quad (30)$$

and we empirically find $\cos$ performs better than $\sin$ as in Eq. (17), probably because the value of $\cos(\theta)$ grows monotonically with $\theta \in [0, \pi]$ but $\sin(\theta)$ does not.

### 3.6 Further Discussions
#### 3.6.1 Learning of problem structure

The inherent working pattern of our proposed neural solver actually learns the underlying structure of graph matching problems. With the connection between graph matching problem and association graph, the original mathematical form of Lawler's QAP transforms into a trackable structure with modern deep learning models. The problem structure is learned by GNN, resulting in a simplified Linear Assignment Problem solved with differentiable Sinkhorn algorithm. In [66] some combinatorial problems over graphs are considered, where problems simplified by GNN are solved greedily. A similar scheme should generalize to other combinatorial problems, by exploiting the representative power of learning problem structures by deep learning.

#### 3.6.2 Matching-aware embedding

The matching-aware Sinkhorn embedding is proposed to add one-to-one matching constraint at shallower embedding layers. Otherwise, the matching constraint is not considered until the output Sinkhorn layer. Early involvement of matching information has been proven effective for both learning-free (RRWM [2] vs SM [7]) and learning based (PCA-GM vs PIA-GM [17]) methods. We show the importance of Sinkhorn embedding by notable improvement in synthetic tests and ablation study on real-world images, additionally further improvement by introducing multi-head Sinkhorn embedding at the cost of increased computation.

#### 3.6.3 Gumbel sampling for optimization problems

As a common post-processing step, the gap between doubly stochastic matrix $\mathbf{S}$ and permutation matrix $\mathbf{X}$ is fulfilled by Hungarian algorithm [43] in a deterministic manner. From the probabilistic point of view, $\mathbf{S}$ represents a distribution on the space of permutation matrices, and our cross-entropy loss minimizes the distance between probability $\mathbf{S}$ and

TABLE 3
Notation of all our proposed methods and their variants.

| model | capability | description |
|---|---|---|
| NGM/-v2 | classic QAP | Neural Graph Matching model introduced in Sec. 3.2 /with enhanced feature in Sec. 3.5. |
| NHGM/-v2 | hypergraph matching | Neural Hyper-Graph Matching model introduced in Sec. 3.3 /with enhanced feature in Sec. 3.5. |
| NMGM/-v2 | multi-graph matching | Neural Multi-Graph Matching model introduced in Sec. 3.4 /with enhanced feature in Sec. 3.5. |
| NGM-V | classic QAP | NGM by replacing Sinkhorn embedding in Eq. (13) with vanilla embedding in Eq. (12). |
| NGM-MH | classic QAP | NGM model with multi-head Sinkhorn embedding (8 multi-head channels). |
| NGM-SF | multi-graph matching | learned NGM model followed with learning-free spectral fusion. |
| NGM-GX | classic QAP | sample X times by Gumbel-Sinkhorn and select the solution with the best objective. |

ground truth distribution $\mathbf{X}^{gt}$. Permutation with the highest probability is selected by Hungarian algorithm.

Such a greedy scheme is empirically successful for matching. However, there might exist better solutions from the distribution, especially when optimizing the objective score of combinatorial problems. Therefore, we switch to Gumbel-Sinkhorn [58] by replacing Eq. (14) with

$$\mathbf{s}_{ia}^{(k)} = \exp\left(\alpha_g(f_c(\mathbf{v}_{ia}^{(k)}) + g)\right) \qquad (31)$$

followed by Sinkhorn algorithm. Note the added $g$ is sampled from standard Gumbel distribution with cumulative distribution function (CDF):

$$G(x) = e^{-e^{-x}} \qquad (32)$$

which models the distribution of extreme values from another distribution. By Eq. (31) sparser doubly-stochastic matrices can be sampled from the original distribution and repeated sampling provides a batch of samples. These sparse matrices are followed by Hungarian discretization, and the objective scores are computed and the best-performing solution is chosen as the final solution. Exploration and speed can be balanced by the number of Gumbel samples.

# 4 EXPERIMENTS

Experiments are conducted on a Linux workstation with Nvidia RTX8000 (48GB) GPU and Intel Xeon W-3175X CPU @ 3.10GHz with 128GB RAM.

We test Lawler's QAP in two settings: i) synthetic point registration, which takes affinity matrix/tensor as input, and ii) QAPLIB with large-scale real-world QAP instances where the network learns to minimize the objective score. We also test our method for a vision application of Lawler's QAP in the sense that CNN features of image keypoints are learned and matched on real images. For keypoint matching, matching accuracy is computed as the percentage of correct matchings among all true matchings.

We also perform hypergraph and multiple graph matching tests to evaluate our NHGM/-v2 and NMGM/-v2. Our PyTorch implementation of NGM/-v2, NHGM/-v2, NMGM/-v2 involves a three-layer GNN, with graph feature channels $l_1 = l_2 = l_3 = 16$. Other hyperparameters are set as $\alpha = \hat{\alpha} = 20, \lambda_2 = 1, \lambda_3 = 1.5$. We set batch size=8. NGM, NHGM, NMGM are trained with SGD and 0.9 Nesterov momentum [67], and the v2 models are trained with Adam optimizer [68]. Detailed learning rate configurations can be found in the following part of this section. Tab. 3 summarizes all our methods and their variants.

## 4.1 Synthetic Experiment for QAP Learning

### 4.1.1 Protocol setting

In the synthetic experiment, sets of random points in the 2D plane are matched by comparison with other competitive learning-free graph matching solvers. For each trial, we generate 10 sets of ground truth points whose coordinates are in the plane $U(0,1) \times U(0,1)$. Synthetic points are distorted by random scaling from $U(1 - \delta_s, 1 + \delta_s)$ and additive random noise $N(0, \sigma_n^2)$. From each ground truth set, 200 graphs are sampled for training and 100 for testing, resulting in totally 2,000 training samples and 1,000 testing samples in one trial. We assume graph structure is unknown to the GM solver, therefore we construct the reference graph by Delaunay triangulation, and the target graph (may contain outliers) is fully connected. Outliers are also randomly sampled from $U(0,1) \times U(0,1)$. By default there are 10 inliers without outlier, with $\delta_s = 0.1, \sigma_n = 0$. We construct the same affinity matrix to formulate Lawler's QAP for all methods.

### 4.1.2 Peer methods

As existing learning methods [17], [25], [31] cannot handle learning with Lawler's QAP with a given affinity matrix, here we first compare learning-free methods: **1) SM** [7] considers graph matching as discovering graph cluster by spectral numerical technique; **2) RRWM** [2] adopts a random-walk view with reweighted jump on graph matching; **3) IPFP** [69] iteratively improves a given solution via integer projection; **4) PSM** [70] improves SM via a probabilistic view; **5) GNCCP** [71] is a convex-concave path-following method for graph matching and **6) BPF** [14] improves path following techniques by branch switching, reaching state-of-the-art performance on graph matching. Additionally, hyper-graph matching algorithm **7) RRWHM** [24] extending powerful RRWM to hyper-graph scenarios is also compared. Second-order affinity is integrated into third-order tensor for RRWHM following [24]. In this experiment, second-order affinity is modeled by $K_{ia,jb} = \exp\left(-(\mathbf{f}_{ij} - \mathbf{f}_{ab})^2/\sigma_2^2\right)$ where $\mathbf{f}_{ij}$ is edge length $E_{ij}$. We empirically set $\sigma_2^2 = 5 \times 10^{-7}$ for all experiments. The third-order affinity model follows Eq. (17) with $\sigma_3 = 0.1$.

For fair comparison of run time, we re-implement the parallelization-friendly SM, RRWM and RRWHM solvers with GPU which can be more scalable than their original single-thread version on CPU. While the other compared methods involve iterative computing and complicated branching, which are not suitable for GPU. Thus the CPU version released by [14] are compared. NGM-V means vanilla NGM without Sinkhorn embedding (see discussion between Eq. (12) and Eq. (13)) and NGM-MH means multi-
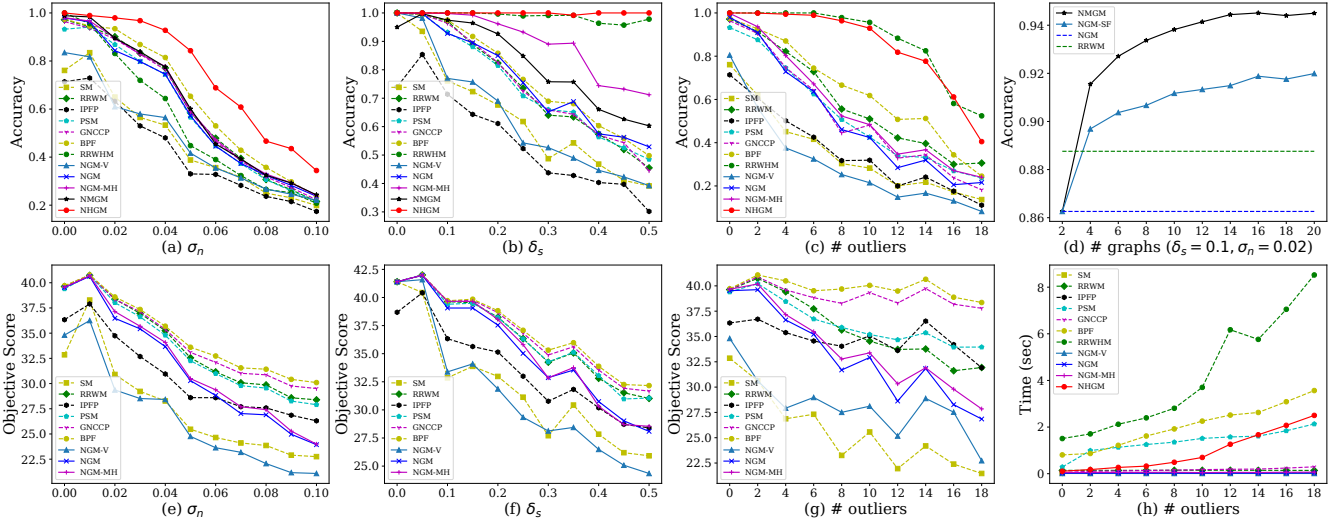
Fig. 6. Synthetic test by varying deformation level $\sigma_n, \delta_s$, number of outliers/graphs. Note it learns a QAP solver based on the given affinity matrix/tensor, without learning any affinity model. This feature is not supported in GMN [25] and PCA-GM [17] thus they cannot be compared.

head Sinkhorn embedding with NGM, by concatenating additional 8 Sinkhorn channels to $\mathbf{m}^{(k)}$ in Eq. (13). The smoothing technique [29] is adopted for multi-matching baseline NGM-SF, where learned NGM is followed with spectral fusion. Since NGM/NHGM/NMGM-v2 share the same solver module with NGM/NHGN/NMGM, they are not compared in synthetic test. The learning rate starts at $10^{-2}$ decays by 10 every 5,000 steps. Multi-graph matching involves 4 graphs by default. The Hungarian algorithm is used as the common discretization step.

### 4.1.3　Result and discussions

Fig. 6(a-c) shows our proposed NGM performs comparatively with state-of-the-art solvers in matching accuracy, and can even surpass under severe random scaling Fig. 6(b). Further improvement in accuracy is achieved via multi-head Sinkhorn embedding model NGM-MH. NMGM gains steadily from NGM by fusing multi-matching information, and in Fig. 6(d) we show the improvement in NMGM by introducing more graphs, and the necessity of learning joint matching as NMGM steadily outperforms NGM-SF whose weights are from two-graph NGM. With the third-order affinity, NHGM shows state-of-the-art robustness to noise, scaling and outliers. Compared to learning-free hyper-graph matching RRWHM [24], our NHGM performs comparatively in the precense of outliers as shown in Fig. 6(c). While it performs more robustly to noises and scaling in Fig. 6(a&b). As shown in Fig. 6(h), NGM, NGM-V and NGM-MH are among the fastest graph matching algorithms, and due to efficient GPU parallelization, NHGM is comparatively fast against PSM [70] and more time-efficient compared to state-of-the-art BPF [14]. In contrast, traditional hypergraph matching algorithms e.g. RRWHM [24] are usually much slower than second-order graph matching and unscalable to larger size of problems.

We report the QAP objective score solved by two-graph matching methods in Fig. 6(e-g), where interestingly our learning-based solvers reach relatively low scores compared to their corresponding accuracy. As have been discussed in

Sec. 1, the QAP objective may be biased under noisy conditions i.e. the optimal solution to QAP may not correspond to true matching. Our solvers learns to ignore the noisy patterns in input affinity matrix. Such a phenomenon becomes more severe in matching real-world images, and the strength of learning-based solvers becomes more significant, as will be shown in Sec. 4.3.1 in details.

The effectiveness of matching-aware Sinkhorn embedding is shown in the accuracy gap between NGM and NGM-V. Further improvement is achieved by multi-head Sinkhorn embedding in NGM-MH, especially with random scaling in Fig. 6(b). As discussed in Sec. 3.6.2, NGM-V without Sinkhorn embedding works in a way similar to SM as the embedding procedure does not consider the assignment constraint, and they also perform closely to each other. On the other hand, by exploiting Sinkhorn embedding, NGM and NGM-MH are conceptually similar to RRWM as all of them try to incorporate the assignment constraint on the fly. Their performances are also very close.

## 4.2　Learning Real-world QAP Instances

### 4.2.1　Experiment setting

Our NGM solves the most general Lawler's QAP, which has a wide range of applications beyond vision. Evaluation on QAPLIB [33] is performed to show the capability of NGM on learning the QAP objective score, which should be minimized in QAPLIB (in contrast, objective score is maximized in graph matching). The QAPLIB contains 134 real-world QAP instances from 15 categories, e.g. planning a hospital facility layout [72]. The problem size is defined as $n_1 = n_2 = n$ from Lawler's QAP in Eq. (1), and ranges from 12 to 256. Results are reported on 133 instances with $12 \le n \le 150$, as the most challenging tai256c is computationally intractable with our testbed (275GB GPU memory is required for intermediate computing). We set the loss function as the objective score of QAP, keeping the

(a) Normalized objective score (lower is better) of our best-performing NGM-G5k against learning-free QAP solvers. Failed instances are plotted at the top of y-axis. The instances are firstly split into two parts based on whether NGM-G5K performs better than Sinkhorn-JA [15] or not, and then sorted by the normalized score of NGM-G5k. NGM-G5k surpasses state-of-the-art learning-free method Sinkhorn-JA [15] on 85 out of 133 instances.



(b) Normalized objective score (lower is better) comparing different sampling settings. More samples in Gumbel-Sinkhorn guarantee a higher probability of finding better solutions, at the cost of increased computation. NGM always picks the permutation matrix with the highest probability with Sinkhorn and Hungarian algorithms, and performs similarly to NGM-G50 and NGM-G500. The order is kept in line with Fig. 7(a).

Fig. 7. Normalized objective score on real-world QAPLIB instances. Only half of the instance labels are shown on x-axis due to limited space.



Fig. 8. Run time (log-scale) against problem size i.e. number of nodes for each graph. All methods except Sinkhorn-JA are implemented and executed on GPUs, as Sinkhorn-JA is challenging to be parallelized.

model architecture unchanged. We formulate a optimization task where the objective score is minimized:

$$L_{obj} = \text{vec}(\mathbf{S})^{\top}\mathbf{K}\text{vec}(\mathbf{S}) \tag{33}$$

where $\mathbf{S}$ is from the output Sinkhorn layer of NGM. Given optimal $L_{obj}$ is reached, the learned $\mathbf{S}$ is a double-stochastically relaxed solution to original QAP. To explore the feasible space, Gumbel-Sinkhorn discussed in Sec. 3.6.3 is adopted during inference.

The naming of QAPLIB instances are based on the following rule: the prefix is the problem category which is the name of the author who proposes the problem, and the number means the size of the problem. If there are multiple problems with the same size, QAPLIB appends a letter starting from a to distinguish.

$$\underline{\text{bur}} \quad - \quad \underline{\text{26}} \quad - \quad \underline{\text{a}}$$
$$\text{author name} \qquad \text{problem size} \qquad \text{index (optional)}$$

We train one network for each set of problems with the same prefix (i.e. problems from the same category), because they usually have common structures (e.g. all bur problems are keyboard layout design problems). And the problem sizes may vary for problems in the same category (e.g. the sizes of esc problems vary from 16 to 128).

We train one model for each category, and report the normalized objective score. In consideration of compact and intuitive illustration, the normalized objective score is computed with the upper bound (primal bound) provided by the up-to-date online benchmark[2] and normalized by the baseline solver spectral matching (SM) [7]:

$$norm\_score = \frac{solved\_score - upper\_bnd}{SM\_score - upper\_bnd} \tag{34}$$

Detailed per-instance scores and timing statics are available in Appendix. Both standard NGM model (NGM) and NGM

2. http://anjos.mgi.polymtl.ca/qaplib/inst.html

TABLE 4
Best-performing occurrence count on QAPLIB among all instances and all tested solvers. Our NGM-G5k surpasses all competing methods on most categories and best performs on 72 out of 133 instances.

| category | bur | chr | els | esc | had | kra | lipa | nug | rou | scr | sko | ste | tai | tho | wil | total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| #instances | 8 | 14 | 1 | 19 | 5 | 3 | 16 | 15 | 3 | 3 | 13 | 3 | 25 | 3 | 2 | 133 |
| SM [7] | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| RRWM [2] | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 |
| Sinkhorn-JA [15] | 0 | 14 | 1 | 1 | 0 | 0 | 15 | 4 | 1 | 0 | 1 | 0 | 7 | 1 | 1 | 46 |
| NGM | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 5 | 0 | 0 | 0 | 3 | 1 | 0 | 0 | 12 |
| NGM-G5 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| NGM-G50 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| NGM-G500 | 1 | 0 | 0 | 3 | 0 | 0 | 0 | 1 | 1 | 0 | 5 | 0 | 1 | 0 | 0 | 12 |
| NGM-G5k | 8 | 1 | 0 | 11 | 5 | 3 | 1 | 10 | 1 | 3 | 7 | 2 | 17 | 2 | 1 | 72 |

TABLE 5
Pearson correlation coefficient $r$ (only for $|r| \geq 0.2$ are shown) between the listed statistics for each problem instance in QAPLIB and the corresponding $gap$ of two methods (see Eq. (35)), as a signal of problem difficulty (bigger gap more difficult). Higher value denotes increased negative effect of the corresponding statistics to the final solution quality. The correlation between statistics and the difference of two methods is also listed on the last row, where higher value denotes NGM-G5k is more sensitive than Sinkhorn-JA and vice versa.

| method | $nz$ | $nz/n^4$ (sparsity) | $\mathbf{K}_{std}/\mathbf{K}_{max}$ | $d_{min}/\bar{\mathbf{K}}$ | $d_{max}/\bar{\mathbf{K}}$ | $d_{std}/\bar{\mathbf{K}}$ | $d_{max}/\bar{d}$ | $d_{std}/\bar{d}$ |
|---|---|---|---|---|---|---|---|---|
| NGM-G5k | 0.130 | **0.631** | **0.212** | **-0.222** | **0.230** | **0.286** | **0.230** | **0.280** |
| Sinkhorn-JA | **0.270** | **0.545** | -0.090 | **-0.220** | **0.291** | **0.355** | **0.291** | **0.475** |
| $gap_{\text{NGM}} - gap_{\text{SJA}}$ | -0.184 | -0.121 | **0.234** | 0.071 | -0.137 | -0.163 | -0.137 | **-0.288** |

with different Gumbel sampling numbers (NGM-GX, X = number of samples) are validated for their performance. The learning rate is initialized at $10^{-4}$ and decays by 10 every 50,000 steps. Batch size is set to 1 and the regularization of Gumbel Sinkhorn $\alpha_g = 1$. Our proposed methods are compared fairly with our GPU implementation of RRWM [2] and SM [7], and results provided in the paper of Sinkhorn-JA [15] (runs on Intel Xeon CPU @ 2.40 GHz). For the problem instances not reported in [15], we assume Sinkhorn-JA fails to reach any feasible solution, as there is no explanation of missing instances in the original paper.

### 4.2.2 Result and analysis

In Fig. 7(a), our method beats RRWM [2] and SM [7] and is comparative and even superior against state-of-the-art Sinkhorn-JA [15]. As there is no learning-based QAP solver, only non-learning methods are compared. The effectiveness of Gumbel sampling discussed in Sec. 3.6.3 is validated in Fig. 7(b), where Gumbel-based NGM-G5k consistently outperforms deterministic NGM, which always picks the permutation with the highest probability by Hungarian algorithm. With decreased sampling number, the performance of Gumbel-based methods gradually degenerates. It suggests that more exploration over sampling space guarantees higher expectations on better solutions.

Further evaluation is given in Tab. 4 and Fig. 8. With learning and Gumbel sampling, our NGM-G5k finds the best solution among 72 out of 133 instances, while state-of-the-art learning-free solver Sinkhorn-JA [15] outperforms on 46 instances so that learning-based solvers e.g. NGM can fit a wider range of problems compared to traditional solvers. More importantly, our best-performing model NGM-G5k is of a magnitude faster than Sinkhorn-JA, and adjusting the sampling number of Gumbel method enables balancing between solution quality and computational demand. Finally, we show the generalization ability among different instances by confusion matrix in Fig. 9, where NGM-G5k



Fig. 9. Generalization test by confusion matrix cross QAP problem instances, where models are learned with instances on y-axis and tested with instances on x-axis. Darker color and lower score correspond to better performance. The tasks are randomly selected from the QAPLIB benchmark. Problem sizes are shown by the numbers in instance names, and our method is insensitive to problem sizes.

is trained and tested on different randomly picked instances. Our model generalizes soundly to unseen instances with different problem sizes. In conclusion, our learning QAP solvers achieve the best accuracy-speed trade-off on QAPLIB and can generalize among different problems.

### 4.2.3 Further discussion

As shown in Tab. 4 and Fig. 7(a), learning-free Sinkhorn-JA [15] and our NGM-G5k performs better on separate categories of QAPLIB, e.g. Sinkhorn-JA performs better on chr and lipa instances while our method is more powerful on bur, esc, nug, scr and tai. Some statistical studies are conducted to discover the relation between model behavior and problem patterns, shedding light for future research on both learning-based and learning-free solvers.

For each problem instance, some statistics are summarized from each instance's affinity matrix: problem size $n$,

**NGM Accuracy (diag: 71.0, all: 41.5)**

| Training \ Testing | bottle | bus | car | cat | chair | dog | sofa | train |
|---|---|---|---|---|---|---|---|---|
| bottle | 74.7 | 36.9 | 35.1 | 18.1 | 23.2 | 18.6 | 20.7 | 49.3 |
| bus | 61.9 | 79.4 | 36.4 | 21.0 | 31.2 | 20.8 | 31.1 | 54.4 |
| car | 61.2 | 42.7 | 79.4 | 20.2 | 29.9 | 20.3 | 37.7 | 61.2 |
| cat | 39.6 | 35.6 | 27.5 | 72.1 | 21.6 | 34.8 | 24.4 | 44.2 |
| chair | 46.0 | 41.9 | 44.7 | 22.0 | 41.7 | 24.3 | 30.0 | 58.6 |
| dog | 42.0 | 36.0 | 30.4 | 52.3 | 24.5 | 58.8 | 27.4 | 50.1 |
| sofa | 66.9 | 51.2 | 41.0 | 22.1 | 30.8 | 23.0 | 72.5 | 63.4 |
| train | 63.0 | 50.1 | 44.1 | 47.0 | 28.7 | 36.5 | 27.3 | 89.9 |

**NGM-v2 Accuracy (diag: 79.5, all: 50.3)**

| Training \ Testing | bottle | bus | car | cat | chair | dog | sofa | train |
|---|---|---|---|---|---|---|---|---|
| bottle | 88.5 | 69.7 | 55.2 | 37.4 | 33.4 | 34.5 | 48.4 | 83.6 |
| bus | 59.0 | 92.4 | 61.0 | 33.3 | 34.2 | 33.0 | 47.3 | 92.1 |
| car | 43.5 | 42.3 | 91.7 | 36.3 | 39.9 | 34.6 | 60.4 | 64.2 |
| cat | 35.8 | 42.8 | 39.3 | 75.4 | 29.5 | 54.5 | 33.0 | 69.0 |
| chair | 54.1 | 59.3 | 59.3 | 35.3 | 55.8 | 35.1 | 57.2 | 73.1 |
| dog | 19.2 | 27.9 | 36.0 | 69.1 | 28.0 | 73.5 | 26.7 | 47.1 |
| sofa | 55.8 | 66.1 | 49.1 | 30.6 | 33.9 | 30.7 | 63.1 | 86.2 |
| train | 26.8 | 65.7 | 44.4 | 23.8 | 28.6 | 24.7 | 40.5 | 95.7 |

**PCA-GM Accuracy (diag: 65.4, all: 52.4)**

| Training \ Testing | bottle | bus | car | cat | chair | dog | sofa | train |
|---|---|---|---|---|---|---|---|---|
| bottle | 81.5 | 66.6 | 43.5 | 41.1 | 28.7 | 29.4 | 27.3 | 71.4 |
| bus | 36.3 | 76.4 | 39.7 | 38.4 | 27.0 | 32.8 | 31.6 | 64.2 |
| car | 46.7 | 65.8 | 69.1 | 46.0 | 29.5 | 36.9 | 46.2 | 62.1 |
| cat | 38.1 | 55.1 | 40.2 | 69.7 | 21.4 | 52.5 | 30.8 | 61.0 |
| chair | 39.5 | 56.2 | 44.0 | 39.5 | 39.9 | 33.8 | 37.8 | 57.8 |
| dog | 44.3 | 63.0 | 45.7 | 63.9 | 29.1 | 65.3 | 38.1 | 65.6 |
| sofa | 49.2 | 59.7 | 47.5 | 44.8 | 28.4 | 39.7 | 40.9 | 71.7 |
| train | 36.8 | 63.3 | 41.0 | 42.3 | 25.3 | 34.3 | 30.5 | 80.3 |

**GMN Accuracy (diag: 40.2, all: 40.5)**

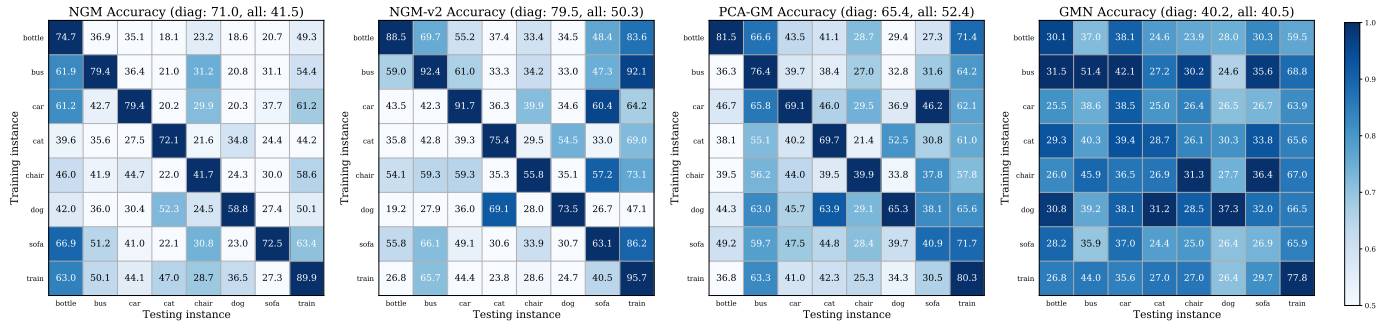| Training \ Testing | bottle | bus | car | cat | chair | dog | sofa | train |
|---|---|---|---|---|---|---|---|---|
| bottle | 30.1 | 37.0 | 38.1 | 24.6 | 23.9 | 28.0 | 30.3 | 59.5 |
| bus | 31.5 | 51.4 | 42.1 | 27.2 | 30.2 | 24.6 | 35.6 | 68.8 |
| car | 25.5 | 38.6 | 38.5 | 25.0 | 26.4 | 26.5 | 26.7 | 63.9 |
| cat | 29.3 | 40.3 | 39.4 | 28.7 | 26.1 | 30.3 | 33.8 | 65.6 |
| chair | 26.0 | 45.9 | 36.5 | 26.9 | 31.3 | 27.7 | 36.4 | 67.0 |
| dog | 30.8 | 39.2 | 38.1 | 31.2 | 28.5 | 37.3 | 32.0 | 66.5 |
| sofa | 28.2 | 35.9 | 37.0 | 24.4 | 25.0 | 26.4 | 26.9 | 65.9 |
| train | 26.8 | 44.0 | 35.6 | 27.0 | 27.0 | 26.4 | 29.7 | 77.8 |

Fig. 10. Generalization study shown as confusion matrix of our proposed NGM and NGM-v2 in line with [17], [25] on Pascal VOC Keypoint. Models are trained on categories on the y-axis and tested on categories on the x-axis. Darker color denotes relatively better performance in the same column, and the averaged accuracy on both diagonal part (learned during training) and all elements (trained+generalized) of confusion matrices are reported in the brackets over confusion matrices. The train/test split follows the main experiment and eight categories are selected randomly.

TABLE 6
Matching accuracy (%) on Pascal VOC Keypoint. Best results are in bold.

| method | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbkie | person | plant | sheep | sofa | train | tv | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GMN [25] | 41.6 | 59.6 | 60.3 | 48.0 | 79.2 | 70.2 | 67.4 | 64.9 | 39.2 | 61.3 | 66.9 | 59.8 | 61.1 | 59.8 | 37.2 | 78.2 | 68.0 | 49.9 | 84.2 | 91.4 | 62.4 |
| PCA-GM [17] | 49.8 | 61.9 | 65.3 | 57.2 | 78.8 | 75.6 | 64.7 | 69.7 | 41.6 | 63.4 | 50.7 | 67.1 | 66.7 | 61.6 | 44.5 | 81.2 | 67.8 | 59.2 | 78.5 | 90.4 | 64.8 |
| IPCA-GM [35] | 53.8 | 66.2 | 67.1 | 61.2 | 80.4 | 75.3 | 72.6 | 72.5 | 44.6 | 65.2 | 54.3 | 67.2 | 67.9 | 64.2 | 47.9 | 84.4 | 70.8 | 64.0 | 83.8 | 90.8 | 67.7 |
| CIE-H [36] | 49.9 | 63.1 | 70.7 | 53.0 | 82.4 | 75.4 | 67.7 | 72.3 | 42.4 | 66.9 | 69.9 | 69.5 | 70.7 | 62.0 | 46.7 | 85.0 | 70.0 | 61.8 | 80.2 | 91.8 | 67.6 |
| LCS [32] | 46.9 | 58.0 | 63.6 | 69.9 | **87.8** | 79.8 | 71.8 | 60.3 | 44.8 | 64.3 | 79.4 | 57.5 | 64.4 | 57.6 | 52.4 | 96.1 | 62.9 | 65.8 | 94.4 | 92.0 | 68.5 |
| BBGM [37] | **61.9** | 71.1 | **79.7** | 79.0 | 87.4 | 94.0 | **89.5** | 80.2 | 56.8 | 79.1 | 64.6 | **78.9** | 76.2 | 75.1 | **65.2** | 98.2 | 77.3 | **77.0** | 94.9 | **93.9** | 79.0 |
| NGM (ours) | 50.1 | 63.5 | 57.9 | 53.4 | 79.8 | 77.1 | 73.6 | 68.2 | 41.1 | 66.4 | 40.8 | 60.3 | 61.9 | 63.5 | 45.6 | 77.1 | 69.3 | 65.5 | 79.2 | 88.2 | 64.1 |
| NHGM (ours) | 52.4 | 62.2 | 58.3 | 55.7 | 78.7 | 77.7 | 74.4 | 70.7 | 42.0 | 64.6 | 53.8 | 61.0 | 61.9 | 60.8 | 46.8 | 79.1 | 66.8 | 55.1 | 80.9 | 88.7 | 64.6 |
| NGM-v2 (ours) | 61.8 | 71.2 | 77.6 | 78.8 | 87.3 | 93.6 | 87.7 | 79.8 | 55.4 | 77.8 | **89.5** | 78.8 | **80.1** | **79.2** | 62.6 | 97.7 | 77.7 | 75.7 | 96.7 | 93.2 | 80.1 |
| NHGM-v2 (ours) | 59.9 | **71.5** | 77.2 | **79.0** | 87.7 | **94.6** | 89.0 | **81.8** | **60.0** | **81.3** | 87.0 | 78.1 | 76.5 | 77.5 | 64.4 | **98.7** | **77.8** | 75.4 | **97.9** | 92.8 | **80.4** |

mean value $\bar{\mathbf{K}}$, minimum value $\mathbf{K}_{min}$, maximum value $\mathbf{K}_{max}$, standard deviation $\mathbf{K}_{std}$, number of zeros $nz$, mean degree (of association graph) $\bar{d}$, minimum degree $d_{min}$, maximum degree $d_{max}$ and standard deviation of degree $d_{std}$. The performance of algorithms is represented by the $gap$ of the solved objective score against upper bound:

$$gap = \frac{solved\_score - upper\_bnd}{solved\_score} \quad (35)$$

It means the percentage of improvement can be made compared to best-known optima (usually solved at extremely high complexity). Pearson correlation coefficients $r$ are computed between each $gap$ and corresponding statistics, additionally some meaningful combinations of the statistics. Items with $|r| \geq 0.2$ are listed in Tab. 5, where positive correlation means a negative effect on solver's performance because a lower $gap$ is better. The correlation between problem statistics and the difference between two methods are also reported.

In the first two columns in Tab. 5, the higher sparsity ($nz/n^4$, proportion of zeros in affinity matrix) makes the problem significantly more challenging for both NGM-G5k and Sinkhorn-JA, and the same conclusion holds for a larger number of zeros $nz$. Then the normalized standard deviation $\mathbf{K}_{std}/\mathbf{K}_{max}$ shows some weak negative effect for NGM-G5k, but little correlation to Sinkhorn-JA. In the last five columns both methods are affected by higher degrees in association graph, while Sinkhorn-JA seems more sensitive to the normalized standard deviation of degrees $d_{std}/\bar{\mathbf{K}}, d_{std}/\bar{d}$. In summary, sparsity is the key challenge for both NGM-G5k (where message passing paths are blocked) and Sinkhorn-JA (where it becomes harder to find tight lower bounds). Furthermore, learning with the association graph can restrain the noisy deviation in degrees. Future improvement may be achieved by designing graph learning models with higher capacity, and designing global communication mechanisms against sparse association graphs.

## 4.3 Real Image for Joint CNN and QAP Learning

Our matching net allows for raw image input, from which a CNN is learned (see Fig. 4). We evaluate semantic keypoint matching on Pascal VOC dataset with Berkeley annotations[3] [73] and Willow ObjectClass dataset [26].

### 4.3.1 Results on Pascal VOC Keypoint dataset

This natural image dataset [73] consists of 20 instance classes with keypoint labels. We follow [17], where image pairs with inlier positions are fed into the model. We consider it a challenging dataset because instance may vary from its scale, pose and illumination, and the number of keypoints in each image varies from 6 to 23. The shallow learning method HARG-SSVM [26] incorporates a fix-sized reference graph, therefore it is inapplicable to our experiment setting where instances from the same category have different inliers. As our multi-matching mechanism in Sec. 3.4 requires the same number of nodes among graphs, our multi-graph model NMGM and NMGM-v2 are not compared either.

In line with the protocol of [17], [25], we filter out those poorly annotated images which are meaningless for matching. Then we get 7,020 training samples and 1,682 testing

3. https://www2.eecs.berkeley.edu/Research/Projects/CS/vision/shape/poselets/voc2011_keypoints_Feb2012.tgz

Fig. 11. Visualization of matching results on 20 Pascal VOC Keypoint categories. Green and red represent correct and incorrect predictions, respectively. As many instances vary a lot in their pose and appearance, this dataset is considered challenging even for deep graph matching.

samples. Instances are cropped around their ground truth bounding boxes and resized to $256 \times 256$ before fed into the network. As discussed in Sec. 3.2.1, we adopt VGG16 backbone [56] and construct the affinity matrix from the same CNN layers: `relu4_2` for node features and `relu5_1` for edge features. The learning rate starts at $10^{-2}$ and decays by 10 every 10,000 steps. For the two input images, one graph is constructed by Delaunay triangulation and the other is fully-connected. $\sigma_3 = 10^{-4}$ is set for third-order affinity of NHGM.

For NGM/NHGM-v2, the feature extractor is replaced by the enhanced backbone with SplineConv and weighted inner-product affinity as discussed in Sec. 3.5, while keeping other modules unchanged. $\sigma_3 = 0.1$ is set for third-order affinity of NHGM-v2. Following [37], we set learning rate $2 \times 10^{-3}$ for VGG16 and $2 \times 10^{-5}$ for other modules.

We compare **GMN** [25], **PCA-GM** [17], **LCS** [32], **BBGM** [37], by which affinity functions are learned for graph matching. It is worth noting that we discover the experiment setting implemented by BBGM [37] is easier than ours, because they filter out keypoints which are out of the bounding box but we do not. Therefore, we reimplement BBGM [37] to fit our experiment setting and the reported BBGM results in this paper are slightly worse than the results in its original paper.

Results in Tab. 6 show that with CNN feature and QAP solver learned jointly, our methods surpass competing methods on most categories, especially best performs in terms of mean accuracy. Specifically, NGM surpasses deep graph matching method PCA-GM, and it is worth noting that PCA-GM incorporates explicit modeling on higher-order and cross-graph affinities, while only second-order affinity is considered in our QAP formulation (and third-order for hypergraph matching). With enhanced feature extractor, our NGM-v2 surpasses state-of-the-art BBGM in terms of accuracy, which is not surprising because NGM-v2 and BBGM share the same feature extractor, and NGM-v2 learns both feature extractor and the graph matching solver, but BBGM only learns the feature extractor. NHGM and NHGM-v bring further improvement with respect to NGM and NGM-v2 by exploiting hypergraph affinities. The GPU memory cost and running speed during training are listed: NGM (4761MB, 13.7 pairs/s); NGM-v2 (5707MB, 14.1 pairs/s); NHGM (9251MB, 10.1 pairs/s); NHGM-v2 (38060MB, 11.0 pairs/s). NGM and NGM-v2 are comparable to PCA-GM [17] (5165MB, 14.4 pairs/s).

The generalization ability is further validated by confusion matrices for NGM and NGM-v2 as shown in Fig. 10. Models are trained only with categories on the x-axis and tested with all categories. The color map is determined

TABLE 7
Controlled experiment by replacing NGM with unlearned ImageNet CNN and RRWM solver. The last row denotes our joint CNN-solver learning method NGM which outperforms.

| CNN | solver | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ImgNet VGG16 [56] | RRWM [2] | 16.1 | 22.3 | 20.8 | 21.8 | 21.3 | 31.0 | 23.2 | 25.4 | 18.6 | 20.5 | 20.6 | 21.7 | 18.8 | 21.9 | 13.5 | 28.6 | 21.7 | 18.3 | 50.5 | 42.8 | 24.0 |
| ImgNet VGG16 [56] | NGM solver (ours) | 30.8 | 42.5 | 44.3 | 33.8 | 39.8 | 52.2 | 49.2 | 53.9 | 27.5 | 42.4 | 29.3 | 49.1 | 45.1 | 45.1 | 24.0 | 48.3 | 49.9 | 29.9 | 70.2 | 73.3 | 44.0 |
| NGM VGG16 (ours) | RRWM [2] | 41.5 | 54.7 | 54.3 | 50.3 | 67.9 | 74.3 | 70.3 | 60.6 | 42.3 | 59.1 | 48.1 | 57.3 | 59.1 | 56.2 | 40.6 | 69.6 | 63.1 | 52.2 | 76.3 | 87.8 | 59.3 |
| NGM VGG16 (ours) | NGM solver (ours) | 50.1 | 63.5 | 57.9 | 53.4 | 79.8 | 77.1 | 73.6 | 68.2 | 41.1 | 66.4 | 40.8 | 60.3 | 61.9 | 63.5 | 45.6 | 77.1 | 69.3 | 65.5 | 79.2 | 88.2 | 64.1 |

TABLE 8
Matching accuracy (%) on Willow ObjectClass dataset. The first 5 rows denote learning-based peer methods. Row 6-9 show our two-graph and hypergraph matching methods. Row 10-11 indicate learning-free multi-graph matching baselines, and HiPPI [74] jointly matches 40, 50, 109, 40, 66 graphs respectively for 5 Willow categories. The last 2 rows denote our two multi-graph matching learning variants.

| model | # graphs | car | duck | face | m-bike | w-bottle | mean |
|---|---|---|---|---|---|---|---|
| GMN [25] | 2 | 67.9 | 76.7 | 99.8 | 69.2 | 83.1 | 79.3 |
| PCA-GM [17] | 2 | 87.6 | 83.6 | 100.0 | 77.6 | 88.4 | 87.4 |
| IPCA-GM [35] | 2 | 90.4 | 88.6 | 100.0 | 83.0 | 88.3 | 90.1 |
| BBGM [37] | 2 | 96.8 | 89.9 | 100.0 | 99.8 | 99.4 | 97.2 |
| LCS [32] | 2 | 91.2 | 86.2 | 100.0 | 99.4 | 97.9 | 94.9 |
| NGM (ours) | 2 | 84.2 | 77.6 | 99.4 | 76.8 | 88.3 | 85.3 |
| NHGM (ours) | 2 | 86.5 | 72.2 | 99.9 | 79.3 | 89.4 | 85.5 |
| NGM-v2 (ours) | 2 | 97.4 | 93.4 | 100.0 | 98.6 | 98.3 | 97.5 |
| NHGM-v2 (ours) | 2 | 97.4 | 93.9 | 100.0 | 98.6 | 98.9 | 97.8 |
| HiPPI [74] | >40 | 74.0 | 88.0 | 100.0 | 84.0 | 95.0 | 88.2 |
| MGM-Floyd [75] | 32 | 85.0 | 79.3 | 100.0 | 84.3 | 93.1 | 88.3 |
| NMGM (ours) | 10 | 78.5 | 92.1 | 100.0 | 78.7 | 94.8 | 88.8 |
| NMGM-v2 (ours) | 10 | **97.6** | **94.5** | **100.0** | **100.0** | 99.0 | **98.2** |

TABLE 9
Ablation study of NHGM-v2 on Pascal VOC Keypoint dataset.

| model | accuracy | relative acc |
|---|---|---|
| baseline NGM | 58.7% | |
| + node affinity | 59.6% | +0.9% |
| + Sinkhorn embedding | 64.1% | +4.5% |
| + SplineConv feature | 74.3% | +10.2% |
| + weighted inner-product affinity | 80.1% | +5.8% |
| + hypergraph affinity | 80.4% | +0.3% |

Our NMGM involves less number of graphs than learning-free methods HiPPI [74] and MGM-Floyd [75], but achieves higher accuracy thanks to end-to-end multi-graph matching learning. Our NMGM-v2 best performs among all methods.

### 4.3.3  Ablation study on Pascal VOC Keypoint dataset

We validate the effectiveness of learning on both CNN and the solver by a controlled experiment. In Tab. 7, the first entry denotes whether the CNN weights are obtained from pretrained ImageNet classifier (ImgNet) [76] or learned graph matching model (NGM); the second entry means either learning-free solver RRWM or learning-based solver NGM is adopted to solve QAP. The necessity of learning on both CNN and the solver is validated, and our joint CNN and solver learning method best performs among them. Learning with CNN unsurprisingly mitigates the gap between classification task and matching task, while it is worth noting the learned solver is nearly twice accurate against RRWM with ImageNet CNN (44.0% vs 24.0%), showing the robustness on our solver side against noises.

Ablation study is performed on our proposed modules in NHGM-v2 on Pascal VOC Keypoint dataset, as shown in Tab. 9. The baseline model is built following NGM introduced in Sec. 3.2, but node affinity is ignored in model input, i.e. $\mathbf{v}_{ia}^{(0)} = 1$ and Sinkhorn embedding (Sec. 3.6.2) is excluded. The effectiveness of node affinity, Sinkhorn embedding, and NGM-v2's SplineConv feature and weighted-inner product affinity are validated by adding these components successively. Finally, we add hypergraph affinity proposed in Sec. 3.5.2 for NHGM-v2. Tab. 9 shows that SplineConv feature contributes significantly to the matching accuracy, and NGM can extend seamlessly to NGM-v2 because our method handles the most general form of QAP.

## 5  CONCLUSION AND OUTLOOK

We have presented a novel neural graph matching network, with three main highlights: i) The first graph matching network directly learning Lawler's QAP which is general with a wide range of applications e.g. on QAPLIB beyond visual matching. This is in contrast to many existing works that can only take separate graphs as input. ii) The first deep network

by the accuracy in current cell normalized by the highest accuracy in its column. As shown in the confusion matrices, both NGM and NGM-v2 own some generalization ability between visually similar categories, e.g. chair and sofa, cat and dog, bus and train. Compared to peer deep graph matching methods [17], [25], NGM and NGM-v2 fit better on the training categories on the diagonal of confusion matrices, while on the other hand NGM and NGM-v2 seem less powerful than PCA-GM when considering all categories, since most categories are irrelevant to the training category.

The matching visualization is given in Fig. 11. The error patterns of NGM-v2, NHGM-v2 are similar but differ from NGM and PCA-GM, because they are based on different feature extractors. As can be seen from the difference between NGM and NGM-v2, better feature extractor can promise better matching accuracy in various categories. And NHGM-v2 improves NGM-v2 by correcting some existing errors, e.g. in "cat", "motorbike", "person" and "tvmonitor".

### 4.3.2  Results on Willow ObjectClass dataset

This natural image dataset covers 5 categories. Each category contains at least 40 images, and all instances in the same class share 10 distinctive image keypoints. We mainly evaluate multi-graph matching learning of NMGM and NMGM-v2 on Willow ObjectClass. Following the protocol in [17], we directly train our methods on the first 20 images and report testing results on the rest. The learning rate starts at $10^{-2}$ and decays by 10 every 500 steps.

The performance of HARG-SSVM [26], GMN [25] and PCA-GM [17] reported in [17] are listed and compared. Two novel multi-graph matching algorithms HiPPI [74] and MGM-Floyd [75] are considered as learning-free baselines. Tab. 8 shows that NGM and NHGM performs comparatively to PCA-GM especially on rigid objects, and NGM-v2 and NHGM-v2 surpass the state-of-the-art BBGM [37].

for hypergraph matching which involves third-order edges. iii) The first network for deep learning of multiple graph matching. Extensive experimental results on synthetic and real-world data show the state-of-the-art performance of our approach. In particular, it shows the notable cost-efficiency advantages against learning-free methods.

In future work, we will explore more scalable approaches for handling large-scale QAP problems. For graph matching, it indicates more graphs with more nodes for matching. For its generality, we will also apply our model to more application areas beyond computer vision.

## ACKNOWLEDGMENTS

## REFERENCES

[1] M. R. Garey and D. S. Johnson, *Computers and Intractability; A Guide to the Theory of NP-Completeness*, 1990.

[2] M. Cho, J. Lee, and K. M. Lee, "Reweighted random walks for graph matching," in *Eur. Conf. Comput. Vis.* Springer, 2010, pp. 492–505.

[3] S. Gold and A. Rangarajan, "A graduated assignment algorithm for graph matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 18, no. 4, pp. 377–388, 1996.

[4] M. A. Fischler and R. C. Bolles, "Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[5] Z. Zhang, "Iterative point matching for registration of free-form curves and surfaces," *Int. J. Comput. Vis.*, vol. 13, no. 2, pp. 119–152, 1994.

[6] E. M. Loiola, N. M. M. de Abreu, P. O. Boaventura-Netto, P. Hahn, and T. Querido, "A survey for the quadratic assignment problem," *Eur. J. Operational Research*, vol. 176, no. 2, pp. 657–690, 2007.

[7] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," in *Int. Conf. Comput. Vis.*, 2005, pp. 1482–1489.

[8] J. Yan, M. Cho, H. Zha, X. Yang, and S. Chu, "Multi-graph matching via affinity optimization with graduated consistency regularization," *Trans. Pattern Anal. Mach. Intell.*, 2016.

[9] C. Edwards, "A branch and bound algorithm for the koopmans-beckmann quadratic assignment problem," in *Combinatorial optimization II.* Springer, 1980, pp. 35–52.

[10] A. P. Punnen and S. N. Kabadi, "A linear time algorithm for the koopmans–beckmann qap linearization and related problems," *Discrete Optimization*, vol. 10, no. 3, pp. 200–209, 2013.

[11] G. Erdoğan and B. Tansel, "A branch-and-cut algorithm for quadratic assignment problems based on linearizations," *Computers & Operations Research*, vol. 34, no. 4, pp. 1085–1106, 2007.

[12] T. Dokeroglu and A. Cosar, "A novel multistart hyper-heuristic algorithm on the grid for the quadratic assignment problem," *Engineering Applications of Artificial Intelligence*, vol. 52, pp. 10–25, 2016.

[13] P. Hahn, T. Grant, and N. Hall, "A branch-and-bound algorithm for the quadratic assignment problem based on the hungarian method," *EJOR*, vol. 108, no. 3, pp. 629 – 640, 1998.

[14] T. Wang, H. Ling, C. Lang, and S. Feng, "Graph matching with adaptive and branching path following," *Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 12, pp. 2853–2867, 2018.

[15] Y. Kushinsky, H. Maron, N. Dym, and Y. Lipman, "Sinkhorn algorithm for lifted assignment problems," *SIAM Journal on Imaging Sciences*, vol. 12, no. 2, pp. 716–735, 2019.

[16] A. Nowak, S. Villar, A. Bandeira, and J. Bruna, "Revised note on learning quadratic assignment with graph neural networks," in *Data Science Workshop*, 2018.

[17] R. Wang, J. Yan, and X. Yang, "Learning combinatorial embedding networks for deep graph matching," in *Int. Conf. Comput. Vis.*, 2019, pp. 3056–3065.

[18] E. L. Lawler, "The quadratic assignment problem," *Management Science*, vol. 9, no. 4, pp. 586–599, 1963.

[19] T. C. Koopmans and M. Beckmann, "Assignment problems and the location of economic activities," *Econometrica*, pp. 53–76, 1957.

[20] M. Chertok and Y. Keller, "Efficient high order matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 12, pp. 2205–2215, 2010.

[21] O. Duchenne, F. Bach, I. Kweon, and J. Ponce, "A tensor-based algorithm for high-order graph matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 12, pp. 2383–2395, 2011.

[22] J. Yan, C. Zhang, H. Zha, W. Liu, X. Yang, and S. M. Chu, "Discrete hyper-graph matching," in *Comput. Vis. Pattern Recog.*, 2015, pp. 1520–1528.

[23] R. Zass and A. Shashua, "Probabilistic graph and hypergraph matching," in *Comput. Vis. Pattern Recog.* IEEE, 2008, pp. 1–8.

[24] J. Lee, M. Cho, and K. M. Lee, "Hyper-graph matching via reweighted random walks," in *Comput. Vis. Pattern Recog.* IEEE, 2011, pp. 1633–1640.

[25] A. Zanfir and C. Sminchisescu, "Deep learning of graph matching," in *Comput. Vis. Pattern Recog.*, 2018, pp. 2684–2693.

[26] M. Cho, K. Alahari, and J. Ponce, "Learning graphs to match," in *Int. Conf. Comput. Vis.*, 2013, pp. 25–32.

[27] J. Yan, J. Wang, H. Zha, X. Yang, and S. Chu, "Consistency-driven alternating optimization for multigraph matching: A unified approach," *Trans. Image Process.*, vol. 24, no. 3, pp. 994–1009, 2015.

[28] Y. Chen, L. Guibas, and Q. Huang, "Near-optimal joint object matching via convex relaxation," in *Int. Conf. Mach. Learn.* PMLR, 2014, pp. 100–108.

[29] D. Pachauri, R. Kondor, and V. Singh, "Solving the multi-way matching problem by permutation synchronization," in *Neural Info. Process. Systems.* Citeseer, 2013, pp. 1860–1868.

[30] Q. Wang, X. Zhou, and K. Daniilidis, "Multi-image semantic matching by mining consistent features," in *Comput. Vis. Pattern Recog.*, 2018, pp. 685–694.

[31] Z. Zhang and W. S. Lee, "Deep graphical feature learning for the feature matching problem," in *Int. Conf. Comput. Vis.*, 2019, pp. 5087–5096.

[32] T. Wang, H. Liu, Y. Li, Y. Jin, X. Hou, and H. Ling, "Learning combinatorial solver for graph matching," in *Comput. Vis. Pattern Recog.*, 2020, pp. 7568–7577.

[33] R. E. Burkard, S. E. Karisch, and F. Rendl, "QAPLIB – a quadratic assignment problem library," *Journal of Global optimization*, vol. 10, no. 4, pp. 391–403, 1997.

[34] E. Maset, F. Arrigoni, and A. Fusiello, "Practical and efficient multi-view matching," in *Int. Conf. Comput. Vis.*, 2017, pp. 4568–4576.

[35] R. Wang, J. Yan, and X. Yang, "Combinatorial learning of robust deep graph matching: an embedding based approach," *IEEE TPAMI*, 2020.

[36] T. Yu, R. Wang, J. Yan, and B. Li, "Learning deep graph matching with channel-independent embedding and hungarian attention," in *Int. Conf. Learn. Rep.*, 2020.

[37] M. Rolínek, P. Swoboda, D. Zietlow, A. Paulus, V. Musil, and G. Martius, "Deep graph matching via blackbox differentiation of combinatorial solvers," in *Eur. Conf. Comput. Vis.* Springer, 2020, pp. 407–424.

[38] P. Swoboda, A. Mokarian, C. Theobalt, F. Bernard *et al.*, "A convex relaxation for multi-graph matching," in *Comput. Vis. Pattern Recog.*, 2019, pp. 11 156–11 165.

[39] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *Trans. Neural Netw.*, vol. 20, no. 1, pp. 61–80, 2008.

[40] J. Yan, S. Yang, and E. R. Hancock, "Learning for graph matching and related combinatorial optimization problems," in *Int. Joint Conf. Artificial Intell.*, 2020.

[41] J. Yan, X.-C. Yin, W. Lin, C. Deng, H. Zha, and X. Yang, "A short survey of recent advances in graph matching," in *Int. Conf. Multimedia Retrieval*, 2016, pp. 167–174.

[42] Y. Bengio, A. Lodi, and A. Prouvost, "Machine learning for combinatorial optimization: a methodological tour d'horizon," *Eur. J. Operational Research*, vol. 290, no. 2, pp. 405–421, 2021.

[43] R. Burkard, M. DellAmico, and S. Martello, *Assignment Problems*. SIAM, 2009.

[44] J. Yan, C. Li, Y. Li, and G. Cao, "Adaptive discrete hypergraph matching," *Trans. Cybern.*, vol. 48, no. 2, pp. 765–779, 2017.

[45] Q. Nguyen, A. Gautier, and M. Hein, "A flexible tensor block coordinate ascent scheme for hypergraph matching," in *Comput. Vis. Pattern Recog.*, 2015, pp. 5270–5278.

[46] J. Yan, Y. Tian, H. Zha, X. Yang, Y. Zhang, and S. M. Chu, "Joint optimization for consistent multiple graph matching," in *Int. Conf. Comput. Vis.*, 2013, pp. 1649–1656.

[47] J. Yan, Y. Li, W. Liu, H. Zha, X. Yang, and S. Chu, "Graduated consistency-regularized optimization for multi-graph matching," in *Eur. Conf. Comput. Vis.* Springer, 2014, pp. 407–422.

[48] T. Yu, J. Yan, W. Liu, and B. Li, "Incremental multi-graph matching via diversity and randomness based graph clustering," in *Eur. Conf. Comput. Vis.*, 2018, pp. 139–154.

[49] L. Torresani, V. Kolmogorov, and C. Rother, "Feature correspondence via graph matching: Models and global optimization," in *Eur. Conf. Comput. Vis.* Springer, 2008, pp. 596–609.

[50] T. Caetano, J. McAuley, L. Cheng, Q. Le, and A. J. Smola, "Learning graph matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 6, pp. 1048–1058, 2009.

[51] M. Leordeanu, R. Sukthankar, and M. Hebert, "Unsupervised learning for graph matching," *Int. J. Comput. Vis.*, vol. 96, no. 1, pp. 28–45, 2012.

[52] M. Leordeanu, A. Zanfir, and C. Sminchisescu, "Semi-supervised learning and optimization for hypergraph matching," in *Int. Conf. Comput. Vis.* IEEE, 2011, pp. 2274–2281.

[53] R. Adams and R. Zemel, "Ranking via sinkhorn propagation," *arXiv:1106.1925*, 2011.

[54] Y. Tian, J. Yan, H. Zhang, Y. Zhang, X. Yang, and H. Zha, "On the convergence of graph matching: Graduated assignment revisited," in *Eur. Conf. Comput. Vis.* Springer, 2012, pp. 821–835.

[55] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *Int. Conf. Learn. Rep.*, 2017.

[56] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Int. Conf. Learn. Rep.*, 2014.

[57] F. Zhou and F. De la Torre, "Factorized graph matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 9, pp. 1774–1789, 2015.

[58] G. Mena, D. Belanger, S. Linderman, and J. Snoek, "Learning latent permutations with gumbel-sinkhorn networks," *Int. Conf. Learn. Rep.*, 2018.

[59] R. Santa Cruz, B. Fernando, A. Cherian, and S. Gould, "Visual permutation learning," *Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 12, pp. 3100–3114, 2018.

[60] F. Serratosa, A. Solé-Ribalta, and X. Cortés, "Automatic learning of edit costs based on interactive and adaptive graph recognition," in *International Workshop on Graph-Based Representations in Pattern Recognition*, 2011, pp. 152–163.

[61] S. Bai, F. Zhang, and P. H. Torr, "Hypergraph convolution and hypergraph attention," *Pattern Recognition*, vol. 110, p. 107637, 2021.

[62] Y. Feng, H. You, Z. Zizhao, R. Ji, and Y. Gao, "Hypergraph neural networks," in *AAAI*, vol. 33, no. 01, 2019, pp. 3558–3565.

[63] X. Zhou, M. Zhu, and K. Daniilidis, "Multi-image matching via fast alternating minimization," in *Int. Conf. Comput. Vis.*, 2015, pp. 4032–4040.

[64] C. Ionescu, O. Vantzos, and C. Sminchisescu, "Matrix backpropagation for deep networks with structured layers," in *Int. Conf. Comput. Vis.*, 2015, pp. 2965–2973.

[65] M. Fey, J. Eric Lenssen, F. Weichert, and H. Müller, "SplineCNN: Fast geometric deep learning with continuous b-spline kernels," in *Comput. Vis. Pattern Recog.*, 2018, pp. 869–877.

[66] E. Khalil, H. Dai, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," in *Neural Info. Process. Systems*, 2017, pp. 6351–6361.

[67] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, "On the importance of initialization and momentum in deep learning," in *Int. Conf. Mach. Learn.* PMLR, 2013, pp. 1139–1147.

[68] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Int. Conf. Learn. Rep.*, Dec 2014.

[69] M. Leordeanu, M. Hebert, and R. Sukthankar, "An integer projected fixed point method for graph matching and map inference," in *Neural Info. Process. Systems.* Citeseer, 2009, pp. 1114–1122.

[70] A. Egozi, Y. Keller, and H. Guterman, "A probabilistic approach to spectral graph matching," *Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 1, pp. 18–27, 2012.

[71] Z.-Y. Liu, H. Qiao, and L. Xu, "An extended path following algorithm for graph-matching problem," *Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1451–1456, 2012.

[72] P. M. Hahn and J. Krarup, "A hospital facility layout problem finally solved," *Journal of Intelligent Manufacturing*, vol. 12, no. 5-6, pp. 487–496, 2001.

[73] L. Bourdev and J. Malik, "Poselets: Body part detectors trained using 3d human pose annotations," in *Int. Conf. Comput. Vis.* IEEE, 2009, pp. 1365–1372.

[74] F. Bernard, J. Thunberg, P. Swoboda, and C. Theobalt, "HiPPI: Higher-order projected power iterations for scalable multi-matching," in *Int. Conf. Comput. Vis.*, 2019, pp. 10 284–10 293.

[75] Z. Jiang, T. Wang, and J. Yan, "Unifying offline and online multi-graph matching via finding shortest paths on supergraph," *Trans. Pattern Anal. Mach. Intell.*, 2020.

[76] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Comput. Vis. Pattern Recog.* IEEE, 2009, pp. 248–255.

**Runzhong Wang** (S'21) is currently a PhD Candidate with Department of Computer Science and Engineering, and AI Institute, Shanghai Jiao Tong University. He obtained B.E. in Electrical Engineering from Shanghai Jiao Tong University. He has published a series of first-authored papers in ICCV 2019, NeurIPS 2020, CVPR 2021, and IEEE TPAMI on machine learning for combinatorial optimization. He serves a reviewer for CVPR 2020/2021, NeurIPS 2020 and AAAI 2021. His research interests include machine learning and combinatorial optimization. He has open-sourced and is maintaining the deep graph matching solver available at: https://github.com/Thinklab-SJTU/PCA-GM which has received more than 400 stars.



**Junchi Yan** (S'10-M'11-SM'21) is currently an Associate Professor with Shanghai Jiao Tong University, Shanghai, China. Before that, he was a Senior Research Staff Member and Principal Scientist with IBM Research – China, where he started his career in April 2011. He obtained the Ph.D. in Electrical Engineering from Shanghai Jiao Tong University. His research interests include machine learning and computer vision. He serves as Area Chair for ICPR 2020, CVPR 2021 and Senior PC for CIKM 2019, IJCAI 2021.



**Xiaokang Yang** (M'00-SM'04-F'19) received the B. S. degree from Xiamen University, in 1994, the M. S. degree from Chinese Academy of Sciences in 1997, and the Ph.D. degree from Shanghai Jiao Tong University in 2000. He is currently a Distinguished Professor of School of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University, Shanghai, China. His research interests include visual signal processing and communication, media analysis and retrieval, and pattern recognition. He serves as an Associate Editor of IEEE Transactions on Multimedia, IEEE Signal Processing Letters. He is a Fellow of IEEE.