

# Graph Edit Distance as a Quadratic Program

Sébastien Bougleux, Benoit Gaüzère, Luc Brun

► **To cite this version:**

Sébastien Bougleux, Benoit Gaüzère, Luc Brun. Graph Edit Distance as a Quadratic Program. ICPR 2016 23rd International Conference on Pattern Recognition, Dec 2016, Cancun, Mexico. IEEE, 2016, Proceedings of ICPR 2016. <<http://www.icpr2016.org/site/>>. <hal-01418937>

**HAL Id: hal-01418937**

**<https://hal.archives-ouvertes.fr/hal-01418937>**

Submitted on 17 Dec 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Graph Edit Distance as a Quadratic Program

Sébastien Bougleux  
Normandie Univ, UNICAEN  
CNRS, GREYC, France

Benoit Gaüzère  
Normandie Univ, INSA Rouen  
LITIS, France

Luc Brun  
Normandie Univ, ENSICAEN  
CNRS, GREYC, France

**Abstract**—The graph edit distance (GED) measures the amount of distortion needed to transform a graph into another graph. Such a distance, developed in the context of error-tolerant graph matching, is one of the most flexible tool used in structural pattern recognition. However, the computation of the exact GED is NP-complete. Hence several suboptimal solutions, such as the ones based on bipartite assignments with edition, have been proposed. In this paper we propose a binary quadratic programming problem whose global minimum corresponds to the exact GED. This problem is interpreted as a quadratic assignment problem (QAP) where some constraints on solutions have been relaxed. This allows to adapt the integer projected fixed point algorithm, initially designed for the QAP, to efficiently compute an approximate GED by finding an interesting local minimum. Experiments show that our method remains quite close to the exact GED for datasets composed of small graphs, while keeping low execution times on datasets composed of larger graphs.

## I. INTRODUCTION

Computing an efficient similarity or dissimilarity measure between graphs is a major problem in structural pattern recognition. The graph edit distance (GED) [1]–[4], developed in the context of error-tolerant graph matching, provides such a measure. It may be understood as the minimal amount of distortion required to transform a graph into another, by a sequence of edit operations applied on nodes and on edges, restricted here to substitutions, insertions and removals. Such a sequence is called an edit path. Depending on the nature of the graphs and on the context, all paths may not have a same importance, in particular when nodes or edges are labeled or attributed. So each possible edit operation  $e$  is penalized by a non-negative cost  $c(e)$ , and the integration of these costs over an edit path  $\gamma$  defines the length  $A(\gamma)$  of this path. An edit path having a minimal length, among all edit paths  $\Gamma(G_1, G_2)$  transforming a graph  $G_1$  into a graph  $G_2$ , defines the GED from  $G_1$  to  $G_2$ :

$$\text{GED}(G_1, G_2) = \min_{\gamma \in \Gamma(G_1, G_2)} \left\{ A(\gamma) \stackrel{\text{def.}}{=} \sum_{e \in \gamma} c(e) \right\} \quad (1)$$

Since computing the GED is NP-complete, it is restricted to rather small graphs. So several approaches have been proposed to approximate the GED efficiently and to process larger graphs. In this paper, graphs are assumed to be simple (no loop nor multiple edge), and each element of the two graphs can be edited only once (no composition of edit operations).

Under these conditions, the GED can be expressed as a binary quadratic program (QP), which is the subject of this paper. Contrary to the graph matching problem, which can be

rewritten as a quadratic assignment problem (QAP) [5], [6], an error-tolerant graph matching can also remove and insert elements. Managing these edit operations is the main difficulty of the quadratic framework. To this, a binary QP is obtained in [7] by considering an edit grid. It is transformed into a binary linear program, which can be solved exactly with classical solvers, but graphs are restricted to be undirected, edges to be unlabeled, and it is also restricted to small graphs since there are  $(n+m)^2$  binary variables, with  $n = |V_1|$  and  $m = |V_2|$ . In [8] a QP with  $nm$  real variables is defined to optimize the cost restricted to node substitutions and associated edge operations. This QP, solved by an interior point algorithm, provides a weight for each possible node substitution. From this fuzzy assignment between nodes, an approximate minimal-cost edit path is then constructed using an heuristic method. Other approximate GED are based on a similar two-step strategy.

In particular, the bipartite graph edit distance (bGED) [9], [10] approximates the GED by replacing the error-tolerant matching between two graphs by an error-tolerant bipartite matching of their nodes, from which an edit path is constructed. Such an error-tolerant set matching can be rewritten as a square linear sum assignment problem (LSAP), with  $(n+m)^2$  binary variables to take into account substitutions, insertions and removals of elements. To approximate the GED, these operations are penalized by the cost of substituting, removing or inserting substructures attached to the nodes [9]–[15]. Since the LSAP can be solved in polynomial time, for instance in  $O((n+m)^3)$  with the well-known Hungarian algorithm [5], [16], the bGED provides an efficient approximation of the GED. This linear framework has been extended to a quadratic one [17], where the GED is expressed as a global solution of a QAP with  $(n+m)^2$  binary variables. An interesting local minimal solution is obtained by relaxing the constraints on the solutions, and by using the integer projected fixed point algorithm (IPFP) [6]. This algorithm iterates: an instance of the LSAP (projection of local continuous solution in the discrete domain), and a line search in the continuous domain to get the next local continuous solution.

While this last quadratic framework improves the approximation of the exact GED, the size of the initial data has been artificially augmented to manage removal and insertion operations more easily. This is also the case in the linear framework above, and in the one proposed by [7]. Recently, the error-tolerant set matching involved in the computation of bGEDs has been formulated as an extension of the LSAP with only  $(n+1)(m+1)$  binary variables [18]. Since a solution can be

computed in  $O(\min\{n, m\}^2 \max\{n, m\})$  time complexity by an adapted Hungarian algorithm, the computation of the bGED is improved, in particular for large graphs.

Based on this last formulation (Sec. II) [18], we propose in this paper a new quadratic expression of the GED with only  $(n+1)(m+1)$  binary variables (Sec. III). This problem is interpreted as a quadratic assignment problem (QAP) where constraints on solutions have been relaxed for removal and insertions operations. This allows to adapt the IPFP algorithm used in [17] to compute efficiently an approximate GED. Experiments (Sec. IV) show that our method remains quite close to the exact GED for datasets composed of small graphs, as previously obtained in [17]. Moreover, contrary to [17], execution times on datasets composed of larger graphs are similar to the ones obtained by bGEDs.

## II. LINEAR SUM ASSIGNMENT PROBLEM WITH EDITION

This section describes the problem of minimal-cost error-tolerant set matching, formalized as an extension of the LSAP, as proposed in [18].

### A. Assignment with edition

A transformation of a set  $V_1$  into a set  $V_2$  can be performed by applying edit operations as follows: Each element  $i \in V_1$  is either substituted by a unique element  $j \in V_2$  (denoted by  $i \rightarrow j$ ), or removed from  $V_1$  ( $i \rightarrow \epsilon$ ). Then each remaining element  $j \in V_2$ , not previously used for a substitution, is inserted (denoted by  $\epsilon \rightarrow j$ ).

Consider the two sets  $V_1^\epsilon = V_1 \cup \{\epsilon\}$  and  $V_2^\epsilon = V_2 \cup \{\epsilon\}$ . Edit operations (substitutions, removals and insertions) defining the transformation can be represented by a mapping  $\varphi: V_1^\epsilon \rightarrow \mathcal{P}(V_2^\epsilon)$  satisfying:

$$\begin{cases} \forall i \in V_1, & |\varphi(i)| = 1 \\ \forall j \in V_2, & |\varphi^{-1}[j]| = 1 \\ & \epsilon \in \varphi(\epsilon) \end{cases} \quad (2)$$

where  $\mathcal{P}(\cdot)$  is the powerset, and  $\varphi^{-1}[j] \stackrel{\text{def.}}{=} \varphi^{-1}[\{j\}]$  denotes the pre-image of any singleton  $\{j\} \in \mathcal{P}(V_2^\epsilon)$  by  $\varphi$ .

Following [18], such a mapping is called an assignment with edition, or  $\epsilon$ -assignment. It may be understood as a bijection on which the bijectivity constraint has been relaxed on the element  $\epsilon$ . Indeed, we implicitly have  $1 \leq |\varphi(\epsilon)| \leq m+1$  and  $1 \leq |\varphi^{-1}[\epsilon]| \leq n+1$ , with  $n = |V_1|$  and  $m = |V_2|$ . Let  $\mathcal{A}_\epsilon(V_1, V_2)$  be the set of  $\epsilon$ -assignments from  $V_1$  to  $V_2$ .

To simplify the forthcoming expressions, consider w.l.o.g. that  $V_1 = \{1, \dots, n\}$ ,  $V_2 = \{1, \dots, m\}$ ,  $V_1^\epsilon = V_1 \cup \{n+1\}$  and  $V_2^\epsilon = V_2 \cup \{m+1\}$ . Any  $\epsilon$ -assignment  $\varphi \in \mathcal{A}_\epsilon(V_1, V_2)$  can be equivalently represented by a  $(n+1) \times (m+1)$  binary matrix

$$\mathbf{X} = \left( \begin{array}{ccc|c} & V_2 & & \epsilon \\ \hline x_{1,1} & \dots & x_{1,m} & x_{1,m+1} \\ \vdots & \ddots & \vdots & \vdots \\ x_{n,1} & \dots & x_{n,m} & x_{n+1,m+1} \\ \hline x_{n+1,1} & \dots & x_{n+1,m} & 1 \end{array} \right) \begin{array}{l} V_1 \\ \epsilon \end{array} \quad (3)$$

such that

$$\begin{cases} x_{i,j} = \delta_{\varphi(i)=\{j\}}, & \forall (i,j) \in V_1 \times V_2 & \text{(sub.)} \\ x_{i,m+1} = \delta_{\varphi(i)=\{\epsilon\}}, & \forall i \in V_1 & \text{(rem.)} \\ x_{n+1,j} = \delta_{\varphi^{-1}[j]=\{\epsilon\}}, & \forall j \in V_2 & \text{(ins.)} \\ x_{n+1,m+1} = 1 & & \end{cases} \quad (4)$$

Due to the constraints on  $\varphi$ , the matrix  $\mathbf{X}$  has a 1 on each of its  $n$  first rows and a 1 on each of its  $m$  first columns:

$$\begin{cases} \sum_{j=1}^{m+1} x_{i,j} = 1, & \forall i = 1, \dots, n \\ \sum_{i=1}^{n+1} x_{i,j} = 1, & \forall j = 1, \dots, m \\ x_{n+1,m+1} = 1 \end{cases} \quad (5)$$

Reciprocally, any matrix  $\mathbf{X} \in \{0, 1\}^{(n+1) \times (m+1)}$  satisfying Eq. 5 represents an  $\epsilon$ -assignment.

### B. Problem formulation

The definition of an  $\epsilon$ -assignment does not rely on the nature of the underlying data. To select a relevant  $\epsilon$ -assignment in  $\mathcal{A}_\epsilon(V_1, V_2)$ , each edit operation  $o$  is penalized by a non-negative cost  $c(o)$ . Then the cost of an  $\epsilon$ -assignment  $\varphi$ , is obtained by summing the costs of its edit operations:

$$\begin{aligned} L(\varphi) &= \sum_{i \in V_1^\epsilon} \sum_{j \in \varphi(i)} c(i \rightarrow j) \\ &= \underbrace{\sum_{\substack{i \in V_1 \\ \varphi(i)=\{j\}}} c(i \rightarrow j)}_{\text{substitutions}} + \underbrace{\sum_{\substack{i \in V_1 \\ \varphi(i)=\{\epsilon\}}} c(i \rightarrow \epsilon)}_{\text{removals}} + \underbrace{\sum_{\substack{j \in V_2 \\ \varphi^{-1}[j]=\{\epsilon\}}} c(\epsilon \rightarrow j)}_{\text{insertions}} \end{aligned} \quad (6)$$

Note that the mapping  $\epsilon \rightarrow \epsilon$  has a zero cost, since it does not represent any edit operation.

The cost of all edit operations, induced by  $\mathcal{A}_\epsilon(V_1, V_2)$ , can be collected into a matrix  $\mathbf{C} \in [0, +\infty)^{(n+1) \times (m+1)}$  such that  $c_{i,j} = c(i \rightarrow j)$  for all  $(i,j) \in V_1^\epsilon \times V_2^\epsilon$ , with  $c_{n+1,m+1} = 0$ . Consider an  $\epsilon$ -assignment  $\varphi$  and its corresponding binary matrix  $\mathbf{X}$  (Eq. 4). Then the cost of  $\varphi$  (Eq. 6) is equal to:

$$L(\mathbf{X}) = \sum_{i=1}^{n+1} \sum_{j=1}^{m+1} c_{i,j} x_{i,j} = \text{vec}(\mathbf{C})^T \text{vec}(\mathbf{X}) \quad (7)$$

where  $\text{vec}(\mathbf{M})$  is the vectorization of the  $p \times q$  matrix  $\mathbf{M}$  obtained by concatenating its rows, *i.e.* a vector of size  $pq$ .

A relevant  $\epsilon$ -assignment is then defined as one having a minimal cost among all  $\epsilon$ -assignments transforming  $V_1$  into  $V_2$ , *i.e.* satisfying

$$\underset{\varphi \in \mathcal{A}_\epsilon(V_1, V_2)}{\text{argmin}} L(\varphi) \quad (8)$$

or equivalently satisfying the following binary linear program:

$$\underset{\mathbf{x}}{\text{argmin}} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{Lx} = \mathbf{1}_{n+m}, \mathbf{x} \in \{0, 1\}^p, x_p = 1 \} \quad (9)$$

where  $p = (n+1)(m+1)$ ,  $\mathbf{c} = \text{vec}(\mathbf{C})$  and the linear system  $\mathbf{Lx} = \mathbf{1}$  is the matrix version of the constraints defining an  $\epsilon$ -assignment (Eq. 2 or Eq. 5). The assignment matrix  $\mathbf{L} \in \{0, 1\}^{(n+m) \times (n+1)(m+1)}$  is defined by:

$$\forall (i,j), \begin{cases} l_{k,(i,j)} = \delta_{k=i}, & \forall k = 1, \dots, n \\ l_{n+k,(i,j)} = \delta_{k=j}, & \forall k = 1, \dots, m \end{cases} \quad (10)$$

Since  $\mathbf{L}$  is totally unimodular, from standard tools in linear programming, the problem defined by Eq. 9 (or Eq. 8) has a binary solution. This problem is close to the classical linear sum assignment problem (LSAP). Compared to the LSAP, elements can also be removed or inserted, *i.e.* the constraints in Eq. 9 are relaxed for one element ( $\epsilon$ ) in each of the two sets. For this reason, the problem defined by Eq. 9 is called the linear sum assignment problem with edition (LSAPE) [18].

The LSAPE can be efficiently solved in polynomial time complexity by adapting algorithms that initially solve the LSAP. An adaptation of the Hungarian algorithm [5], [16] is presented in [18]. It computes a solution to the LSAPE in  $O(\min\{n, m\}^2 \max\{n, m\})$  time complexity and in  $O(nm)$  space complexity. This improves complexities obtained by the classical approach used to transform a set into another with edit operations [9], [10]. In this last approach, the problem is formalized as a LSAP with a specific  $(n+m) \times (m+n)$  cost matrix. It can be solved by the classical Hungarian algorithm in  $O((n+m)^3)$  time complexity and in  $O((n+m)^2)$  space complexity.

### III. A QUADRATIC EXPRESSION OF THE GED

In this section we extend the linear sum assignment problem with edition to a quadratic one. This simplifies the quadratic expression of the GED proposed in [17].

#### A. Simultaneous assignment with edition

Let  $G_1 = (V_1, E_1)$  and  $G_2 = (V_2, E_2)$  be two simple graphs, directed or not. A transformation of  $G_1$  into  $G_2$ , with edit operations applied on nodes and on edges, can be fully represented by an  $\epsilon$ -assignment  $\varphi: V_1^\epsilon \rightarrow \mathcal{P}(V_2^\epsilon)$  transforming  $V_1$  into  $V_2$ . Then the set  $E_1$  is transformed into the set  $E_2$  as follows:

- 1) Any edge  $(i, j) \in E_1$  is mapped by  $\varphi$  to a pair  $(k, l) \in V_2^\epsilon \times V_2^\epsilon$  with  $\varphi(i) = \{k\}$  and  $\varphi(j) = \{l\}$ . The edge  $(i, j)$  is either
  - a) substituted by  $(k, l)$  if  $(k, l) \in E_2$ , or
  - b) removed if  $(k, l) \notin E_2$ , in particular if  $i$  or  $j$  is removed by  $\varphi$  ( $k = \epsilon$  or  $l = \epsilon$ ).
- 2) Any edge  $(k, l) \in E_2$  is inserted into  $E_1$  if it has not been used by a substitution in Step 1, *i.e.* if there is a pair  $(i, j) \in V_1^\epsilon \times V_1^\epsilon$  such that  $k \in \varphi(i)$ ,  $l \in \varphi(j)$  and  $(i, j) \notin E_1$ . In particular, this occurs when  $k$  or  $l$  is inserted by  $\varphi$  into  $V_1$  ( $i = \epsilon$  or  $j = \epsilon$ ).

The sequence of edit operations transforming  $E_1$  into  $E_2$  is thus represented by a simultaneous  $\epsilon$ -assignment of the nodes. Its cost is defined by summing the cost of its edit operations:

$$\begin{aligned}
Q(\varphi) = & \sum_{\substack{(i,j) \in E_1 \\ \text{with } \{k\}=\varphi(i), \{l\}=\varphi(j), (k,l) \in E_2}} c((i,j) \rightarrow (k,l)) \quad (\text{sub.}) \\
& + \sum_{\substack{(i,j) \in E_1 \\ \text{with } \{k\}=\varphi(i), \{l\}=\varphi(j), (k,l) \notin E_2}} c((i,j) \rightarrow \epsilon) \quad (\text{rem.}) \\
& + \sum_{\substack{(k,l) \in E_2 \\ \text{with } \{i\}=\varphi^{-1}[k], \{j\}=\varphi^{-1}[l], (i,j) \notin E_1}} c(\epsilon \rightarrow (k,l)) \quad (\text{ins.})
\end{aligned} \tag{11}$$

TABLE I

EDIT OPERATION INDUCED BY THE SIMULTANEOUS ASSIGNMENT OF  $i, j \in V_1^\epsilon$  TO  $k, l \in V_2^\epsilon$  RESPECTIVELY

$(i, j)$	$(k, l)$	edit operation	cost $d_{ik,jl}$
$\in E_1$	$\in E_2$	substitution of $(i, j)$ by $(k, l)$	$c((i, j) \rightarrow (k, l))$
$\in E_1$	$\notin E_2$	removal of $(i, j)$	$c((i, j) \rightarrow \epsilon)$
$\notin E_1$	$\in E_2$	insertion of $(k, l)$ into $E_1$	$c(\epsilon \rightarrow (k, l))$
$\notin E_1$	$\notin E_2$	do nothing	0

More generally, the cost of simultaneously assigning  $i \in V_1^\epsilon$  to  $k \in V_2^\epsilon$ , and  $j \in V_1^\epsilon$  to  $l \in V_2^\epsilon$ , can be defined by:

$$\begin{aligned}
d_{ik,jl} = & c((i, j) \rightarrow (k, l)) \delta_{(i,j) \in E_1} \delta_{(k,l) \in E_2} \quad (\text{sub.}) \\
& + c((i, j) \rightarrow \epsilon) \delta_{(i,j) \in E_1} (1 - \delta_{(k,l) \in E_2}) \quad (\text{rem.}) \\
& + c(\epsilon \rightarrow (k, l)) (1 - \delta_{(i,j) \in E_1}) \delta_{(k,l) \in E_2} \quad (\text{ins.})
\end{aligned} \tag{12}$$

where  $\delta_r = 1$  if the relation  $r$  is true or  $\delta_r = 0$  else. The edit operation associated to the simultaneous node assignment is reported in Table I. Observe that  $d_{ik,jl} = 0$  when both  $(i, j)$  and  $(k, l)$  are not edges of  $G_1$  and  $G_2$  respectively. Indeed in this case the simultaneous assignment does not correspond to an edit operation applied on edges. Then the cost  $Q(\varphi)$  defined by Eq. 11 can be written more compactly as:

$$\begin{aligned}
Q(\varphi) = & g \sum_{i \in V_1^\epsilon} \sum_{k \in \varphi(i)} \sum_{j \in V_1^\epsilon} \sum_{l \in \varphi(j)} d_{ik,jl} \\
& \text{with } g = \begin{cases} \frac{1}{2} & \text{if } G_1 \text{ and } G_2 \text{ are undirected} \\ 1 & \text{else} \end{cases}
\end{aligned} \tag{13}$$

When graphs are both undirected, the cost  $d_{ik,jl}$  is symmetric ( $d_{ik,jl} = d_{jl,ik}$ ). So a factor  $\frac{1}{2}$  is introduced to count the cost of each edit operation only once.

#### B. Quadratic assignment problem with edition and GED

An optimal transformation of  $G_1$  into  $G_2$  consists in computing an  $\epsilon$ -assignment having a minimal cost that penalizes operations both on nodes and edges, among all  $\epsilon$ -assignments from  $V_1$  onto  $V_2$ . In other terms, the cost of this optimal transformation is equal to the GED from  $G_1$  to  $G_2$  [17]:

$$\text{GED}(G_1, G_2) = \min_{\varphi \in \mathcal{A}_\epsilon(V_1, V_2)} \left\{ A(\varphi) \stackrel{\text{def.}}{=} Q(\varphi) + L(\varphi) \right\} \tag{14}$$

This optimization problem is close to the quadratic assignment problem (QAP), see [5] for more details on QAP.

To see this, consider an  $\epsilon$ -assignment  $\varphi \in \mathcal{A}_\epsilon(V_1, V_2)$  and the associated binary vector  $\mathbf{x} \in \{0, 1\}^p$ , with  $x_p = 1$  and  $p = (n+1)(m+1)$ . Nodes  $i$  and  $j$  of  $V_1^\epsilon$  are simultaneously assigned by  $\varphi$  to  $k$  and  $l$  of  $V_2^\epsilon$ , respectively, iff  $x_{ik} = x_{jl} = 1$ .

The quadratic cost  $Q(\varphi)$  can then be rewritten as

$$Q(\mathbf{x}) = g \sum_{i=1}^{n+1} \sum_{k=1}^{m+1} \sum_{j=1}^{n+1} \sum_{l=1}^{m+1} d_{ik,jl} x_{ik} x_{jl} = g \mathbf{x}^T \mathbf{D} \mathbf{x} \tag{15}$$

where the matrix  $\mathbf{D} \in [0, +\infty)^{p \times p}$ , of general term  $d_{ik,jl}$  (Eq. 12), encodes the cost of all possible simultaneous assignments. Compared to Eq. 13, the sum in Eq. 15 is also

carried on simultaneous assignments which are not induced by the  $\epsilon$ -assignment  $\varphi$ , *i.e.*  $x_{ik} = 0$  or  $x_{jl} = 0$ . Since they do not contribute to the sum, we have  $Q(\mathbf{x}) = Q(\varphi)$ . Due to the constraints on  $\epsilon$ -assignments, when  $i, j \in V_1$  and  $k, l \in V_2$  remark that several simultaneous assignments never occur, *i.e.*  $x_{ik}x_{il} = 0$  with  $k \neq l$  or  $x_{ik}x_{jk} = 0$  with  $i \neq j$  (a node cannot be assigned twice). In these cases, the value of  $d_{ik,jl}$  is not well-defined by Eq. 12 and should be equal to 0. Nevertheless  $d_{ik,jl}x_{ik}x_{jl} = 0$  is always satisfied according to the previous remark, and so the value of the cost is not so relevant. Also, since graphs are simple, all diagonal elements are null ( $d_{ik,ik} = 0$ ).

The cost of transforming  $G_1$  into  $G_2$  can then be written as

$$A(\mathbf{x}) = g \mathbf{x}^T \mathbf{D} \mathbf{x} + \mathbf{c}^T \mathbf{x} \quad (16)$$

where  $\mathbf{c} \in [0, +\infty)^p$  is the edit cost vector defined in Eq. 9. Since  $\mathbf{x}^T \mathbf{D} \mathbf{x} = \frac{1}{2} \mathbf{x}^T (\mathbf{D} + \mathbf{D}^T) \mathbf{x}$ , we have

$$A(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \Delta \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

with  $\Delta = \begin{cases} \mathbf{D} & \text{if } G_1 \text{ and } G_2 \text{ are undirected} \\ \mathbf{D} + \mathbf{D}^T & \text{else} \end{cases}$  (17)

The cost matrix  $\Delta$  is always symmetric, even if at least one of the two graphs is directed. The linear part of the cost can also be included in the quadratic part, leading to

$$A(\mathbf{x}) = \mathbf{x}^T \widehat{\Delta} \mathbf{x}, \text{ with } \widehat{\Delta} = \frac{1}{2} \Delta + \text{diag}(\mathbf{c}) \quad (18)$$

Then the GED is rewritten as the following quadratic program:

$$\text{GED}(G_1, G_2) = \min \left\{ \mathbf{x}^T \widehat{\Delta} \mathbf{x} \mid \mathbf{L} \mathbf{x} = \mathbf{1}_{n+m}, \mathbf{x} \in \{0, 1\}^p \right\} \quad (19)$$

We call it the quadratic assignment problem with edition (QAPE). The QAPE differs from the classical QAP by the constraint matrix  $\mathbf{L}$ . A QAP finds a bijection or an injection inducing a simultaneous assignment, as for the quadratic expression of the GED proposed in [17]. This expression describes a QAP involving assignments represented by vectors of size  $(n+m)^2$ . This is simplified by Eq. 19 by involving vectors of size  $(n+1)(m+1)$  only.

### C. Relaxation, approximation, algorithm

While exact solutions to the QAPE could be computed by adapting algorithms that initially solve the QAP [5], this would not be competitive with bipartite GEDs in terms of computational time. So we propose to approximate the GED by adapting the integer projected fixed point (IPFP) algorithm [6], designed to approximate the solution to the QAP, and it is also used in [17] to approximate the GED. As several other algorithms, the IPFP algorithm tries to find a solution to the relaxed quadratic problem given here by

$$\text{argmin} \left\{ \mathbf{x}^T \widehat{\Delta} \mathbf{x} \mid \mathbf{L} \mathbf{x} = \mathbf{1}_{n+m}, \mathbf{x} \in [0, 1]^p \right\} \quad (20)$$

where the space of solutions is now continuous.

Given an initial candidate solution  $\mathbf{x}$ , the algorithm iterates the three following steps

$$\mathbf{b}^* \leftarrow \text{argmin} \left\{ (\mathbf{x}^T \widehat{\Delta}) \mathbf{b} \mid \mathbf{L} \mathbf{b} = \mathbf{1}, \mathbf{b} \in \{0, 1\}^p, b_p = 1 \right\} \quad (21)$$

$$\alpha^* \leftarrow \text{argmin}_{\alpha \in [0, 1]} A(\mathbf{x} + \alpha(\mathbf{b}^* - \mathbf{x})) \quad (22)$$

$$\mathbf{x} \leftarrow \mathbf{x} + \alpha^*(\mathbf{b}^* - \mathbf{x}) \quad (23)$$

until convergence (when  $\mathbf{x}$  becomes discrete) or after a given number of iterations.

The first step (Eq. 21) consists in computing a linear approximation of  $A$  around the current solution  $\mathbf{x}$  in the discrete domain. This can be realized by solving a LSAP with  $\mathbf{x}^T \widehat{\Delta}$  as an edit cost matrix (Sec. II). Indeed, the first-order Taylor expansion of  $A$  around  $\mathbf{x}$  is given by

$$A(\mathbf{y}) \approx A(\mathbf{x}) + (\mathbf{x}^T \widehat{\Delta})(\mathbf{y} - \mathbf{x}), \mathbf{y} \in [0, 1]^p \quad (24)$$

so minimizing  $A(\mathbf{x})$  locally is approximatively equivalent to minimizing  $\mathbf{x}^T \widehat{\Delta} \mathbf{y}$  with  $\mathbf{x}$  fixed. By standard tools in linear programming, this linear program has a same binary solution  $\mathbf{b}^*$  as the problem defined by Eq. 21, here a LSAP. It can be solved in  $O(\min\{n, m\}^2 \max\{n, m\})$  time complexity and in  $O(nm)$  space complexity (Sec. II-B). Compared to the IPFP algorithm proposed in [17], where Eq. 21 corresponds to a square LSAP, computational time is reduced at each iteration. This confirmed experimentally (Section IV).

The second step consists in computing the local minimum of  $A$ , by a line search, between the current solution  $\mathbf{x}$  and the candidate binary solution  $\mathbf{b}^*$  included. As detailed in [17] this can be solved analytically, here in  $O(mn)$ . The vector realizing this minimum is the new candidate solution.

The iterative process generally converges to a local minimum of the relaxed problem defined by Eq. 20. When  $\mathbf{x}$  is continuous, an additional projection (Eq. 21) is performed to get the final binary vector. The quality of the approximation depends mainly on the initialization. A random vector satisfying the constraints could be used. Instead, it seems to be more natural to compute a bipartite GED (bGED), *i.e.* to solve a LSAP [18], and consider the corresponding optimal  $\epsilon$ -assignment as initial vector in IPFP algorithm. In this way, results obtained by a bGED should be improved. This is analyzed experimentally in the following section.

## IV. EXPERIMENTS

### A. Chemical datasets

Table II shows results obtained by different state of the art methods to compute an estimation of the GED. These methods have been tested on four chemoinformatics datasets<sup>1</sup>, in the same way as in previous papers using these datasets [12], [13]. Each dataset is composed of molecular graphs. Costs have been set to 1 for any substitution operation, and 3 for insertions and deletions. It is important to notice that each graph adjacency matrix is randomly permuted in order

<sup>1</sup>Datasets are available at <https://iapr-tc15.greyc.fr/links.html>

TABLE II  
ACCURACY AND COMPUTATIONAL SCORES.  $d$  IS THE AVERAGE EDIT DISTANCE,  $e$  THE AVERAGE ERROR AND  $t$  THE AVERAGE COMPUTATIONAL TIME.

	Algorithm	Alkane			Acyclic			MAO		PAH	
		$d$	$e$	$t$	$d$	$e$	$t$	$d$	$t$	$d$	$t$
1	$A^*$	15		1.29	17		6.02				
2	[9], [10]	35	18	$\simeq 10^{-3}$	35	18	$\simeq 10^{-3}$	105	$\simeq 10^{-3}$	138	$\simeq 10^{-3}$
3	[12]	33	18	$\simeq 10^{-3}$	31	14	$\simeq 10^{-2}$	49	$\simeq 10^{-2}$	120	$\simeq 10^{-2}$
4	[13]	26	11	2.27	28	9	0.73	44	6.16	129	2.01
5	QAP	20.2	4.7	0.002	20.8	3.4	0.004	33.9	0.019	53.9	0.038
5	(this paper) QAPE	20.1	4.6	$\simeq 10^{-3}$	20.8	3.5	0.002	33.9	0.005	53.6	0.010

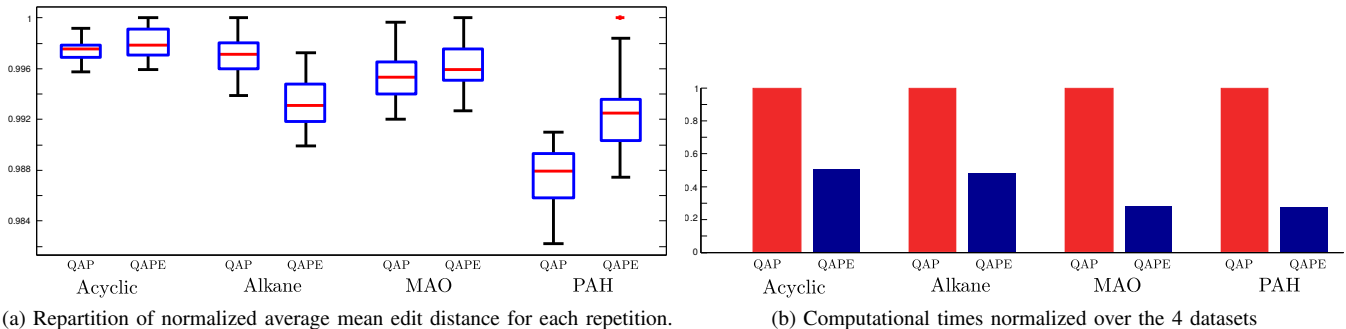


Fig. 1. Comparison of IPFP algorithm for the QAP and for the QAPE.

to avoid some alignment bias which may provide optimistic approximations of the GED. The tested methods are:

- Exact GED based on  $A^*$  algorithm (line 1), restricted to very small graphs (not adapted to MAO and PAH).
- Bipartite GED based on the Hungarian algorithm adapted to the LSAPE [18], with the substructure attached to each node defined by: incident edges and adjacent nodes [9], [10], a bag of random walks up to 3 edges (line 2) [12], a subgraph up to a radius 3 (line 3) [13].
- Approximate GED formalized as a QAP (line 4) [17] and as a QAPE (line 5), both solved by the IPFP algorithm with the same initialization provided by [12], since initialization influences the approximation of the GED.

Observe that these two last approaches obtain the best approximation results. Considering Alkane and Acyclic datasets, this is trivially inferred from average errors. For the two remaining datasets, we observe lower average edit distances which may correspond to better approximations since the exact GED is always over-estimated. In addition, we can note that computational times of quadratic approaches remain globally acceptable, especially compared to the bGED proposed in [13] where most of the time is spent to compute the costs.

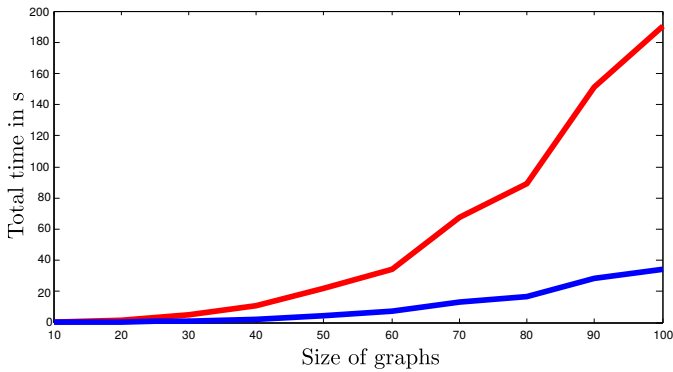
Moreover, quadratic approaches seem to have similar accuracies, while the approach presented in this paper improves clearly the computational time. This is analyzed through a second experiment consisting in 30 repetitions of each approach, for each dataset. As before, graph adjacency matrices are randomly permuted in order to avoid some alignment bias. Figure 1a shows the boxplot of the normalized average approximate GEDs. For each dataset, these average values are

normalized by the maximal mean approximate GED given by  $\max\{\max_i\{\bar{d}_i^{\text{QAPE}}\}, \max_i\{\bar{d}_i^{\text{QAP}}\}\}$  where  $\bar{d}_i$  is the mean approximate GED obtained for the repetition  $i \in \{1, \dots, 30\}$ . Considering Acyclic and MAO datasets, the distributions are statistically similar according to a paired Student's test with a  $p$ -value equals to 0.01. Considering Alkane and PAH datasets, the difference is statistically significant but not so important for the approximation of GEDs (about 0.4%). From a computational time point of view, Figure 1b shows the normalized mean times according to a same initialization. We can clearly see that the proposed approach is faster than previous one, with a ratio around 2 for smaller graphs and 4 for larger graphs.

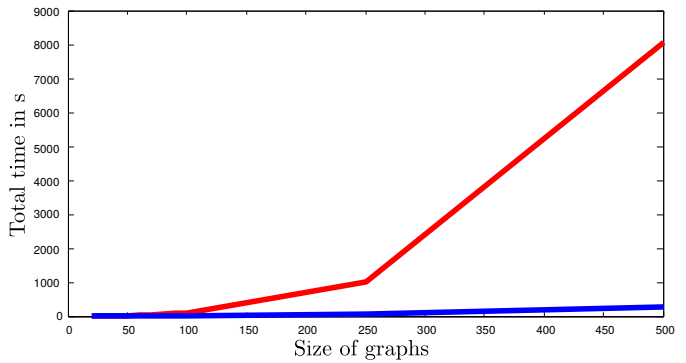
### B. Synthetic datasets

To explore the behavior of the quadratic approaches when the number of nodes increases, we generate synthetic datasets having the same characteristics as MAO dataset but generalized to different graph sizes. These characteristics include node's and edge's labels distribution and ratio between number of edges and number of nodes. For a given number of nodes, each new synthetic graph is created according to these characteristics. This allows to retrieve characteristics of common molecular graphs. For each tested size, up to 100 nodes, a dataset composed of 100 synthetic graphs has been synthesized. Each of this graph, called source graph, is then associated to a target graph which has been obtained by removing one node and by substituting another one from the associated source graph. The exact GED between a source and a target is thus known by construction and is around 10.

Figure 2a shows the behavior of the quadratic approaches when the size increases. As we can see, the proposed approach



(a) Both graphs have the same size (up to one node)



(b) 1st graph has 20 nodes, 2nd second one has 20 to 500 nodes

Fig. 2. Total computational time in seconds taken by IPFP algorithm for the QAP (upper red line) and for the QAPE (bottom blue).

based on a QAPE scales better in terms of computational times. Indeed, for graphs having 100 nodes, the computational ratio is over 5. Moreover, the average error is the same for both methods and remains quite stable on the set of sizes. It is worth noticing that the most important part of time is dedicated here to the initialization, realized as in the previous section. Figure 2b shows a similar result but with a fixed size for the first graph (20 nodes) and a varying size for the second one (from 20 to 500 nodes). The gain in computational times increases as the size of graphs also increases. Considering a second graph with 500 nodes, this gain corresponds to a computational time reduced by a factor 29. This is explained by the reduction of variable sizes from  $(n+m)^2$  to  $(m+1)(n+1)$ , and by the use of a LSAP instead of a LSAP, at each iteration of IPFP algorithm (Sec. III-C).

## V. CONCLUSION

Based on the notion of assignment with edition, this paper presents a binary quadratic programming problem whose global optimum is equal to the GED when graphs are simple, and when edit operations are restricted to substitutions, removals and insertions. The proposed quadratic program reformulates a previous expression of the GED as a QAP, by reducing the number of binary variables to  $(n+1)(m+1)$  and by relaxing the constraints associated to removals and insertions. This allows us to adapt an iterative algorithm, that initially finds an approximate solution to the QAP, to approximate the GED and reduce both space and time complexities. This is confirmed experimentally. Experiments also show that computational times are close to those of bGEDs, while providing distances closer to exact GEDs, on tested datasets.

## REFERENCES

- [1] W.-H. Tsai and K.-S. Fu, "Error-correcting isomorphisms of attributed relational graphs for pattern analysis," *IEEE Trans. on Systems, Man and Cybernetics*, vol. 9, no. 12, pp. 757–768, 1979.
- [2] H. Bunke and G. Allermann, "Inexact graph matching for structural pattern recognition," *Pattern Recognition Letters*, vol. 1, no. 4, pp. 245–253, 1983.
- [3] A. Sanfeliu and K.-S. Fu, "A distance measure between attributed relational graphs for pattern recognition," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 13, no. 3, pp. 353–362, 1983.
- [4] H. Bunke, "Error correcting graph matching: on the influence of the underlying cost function," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 21, no. 9, pp. 917–922, 1999.
- [5] R. Burkard, M. Dell'Amico, and S. Martello, *Assignment Problems*. SIAM, 2009.
- [6] M. Leordeanu, M. Hebert, and R. Sukthankar, "An integer projected fixed point method for graph matching and map inference," in *Advances in Neural Information Processing Systems*, 2009, vol. 22, pp. 1114–1122.
- [7] D. Justice and A. Hero, "A binary linear programming formulation of the graph edit distance," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 8, pp. 1200–1214, 2006.
- [8] M. Neuhaus and H. Bunke, "A quadratic programming approach to the graph edit distance problem," in *Graph-Based Representations in Pattern Recognition*, ser. LNCS, 2007, vol. 4538, pp. 92–102.
- [9] K. Riesen, M. Neuhaus, and H. Bunke, "Bipartite graph matching for computing the edit distance of graphs," in *Graph-Based Representations in Pattern Recognition*, ser. LNCS, 2007, vol. 4538, pp. 1–12.
- [10] K. Riesen and H. Bunke, "Approximate graph edit distance computation by means of bipartite graph matching," *Image and Vision Comp.*, vol. 27, pp. 950–959, 2009.
- [11] Z. Zeng, A. K. H. Tung, J. Wang, J. Feng, and L. Zhou, "Comparing stars: On approximating graph edit distance," *Proceedings of the VLDB Endowment*, vol. 2, no. 1, pp. 25–36, 2009.
- [12] B. Gaüzère, S. Bougleux, K. Riesen, and L. Brun, "Approximate Graph Edit Distance Guided by Bipartite Matching of Bags of Walks," in *Structural, Syntactic and Statistical Pattern Recognition*, ser. LNCS, 2014, vol. 8621, pp. 73–82.
- [13] V. Carletti, B. Gaüzère, L. Brun, and M. Vento, "Approximate graph edit distance computation combining bipartite matching and exact neighborhood substructure distance," in *Graph-Based Representations in Pattern Recognition*, ser. LNCS, 2015, vol. 9069, pp. 168–177.
- [14] K. Riesen and H. Bunke, "Improving bipartite graph edit distance approximation using various search strategies," *Pattern Recognition*, vol. 28, no. 4, pp. 1349–1363, 2015.
- [15] K. Riesen, A. Fischer, and H. Bunke, "Estimating graph edit distance using lower and upper bounds of bipartite approximations," *Int. J. Patt. Recogn. Artif. Intell.*, vol. 29, no. 2, 2015.
- [16] E. Lawler, *Combinatorial Optimization: Networks and Matroids*. New York: Holt, Rinehart and Winston, 1976.
- [17] S. Bougleux, L. Brun, V. Carletti, P. Foggia, B. Gaüzère, and M. Vento, "A quadratic assignment formulation of the graph edit distance," Normandie Univ, NormaSTIC FR 3638, France, Tech. Rep., 2015. [Online]. Available: <http://arxiv.org/abs/1512.07494>
- [18] S. Bougleux and L. Brun, "Linear sum assignment with edition," Normandie Univ, GREYC UMR 6072, Tech. Rep., 2016.