

Evaluation of Graph Matching Measures for Documents Retrieval*

Salim Jouili*, Salvatore Tabbone* and Ernest Valveny⁺

* *LORIA UMR-7503, University of Nancy 2,
BP 239, 54506 Vandoeuvre-lès-Nancy Cedex, France
E-mail: {salim.jouili, tabbone}@loria.fr*

⁺ *Centre de Visió per Computador, Dep. Ciències de la Computació
Universitat Autònoma de Barcelona, Edifici O, Campus UAB, 08193 Bellaterra, Spain
E-mail: ernest@cvc.uab.cat*

Abstract

In this paper we evaluate four graph distance measures. The analysis is performed for document retrieval tasks. For this aim, different kind of documents are used which include line drawings (symbols), ancient documents (ornamental letters), shapes and trademark-logos. The experimental results show that the performance of each graph distance measure depends on the kind of data and the graph representation technique.

Keywords: Graph matching, Graph retrieval, Structural representation, Performance evaluation.

1 Introduction

In document retrieval applications, it is necessary to define some description of the document based on a set of features. These descriptions are then used to search and to determine which documents satisfy the query selection criteria. The effectiveness of a document retrieval system ultimately depends on the type of representation used to describe a document. In pattern recognition, the document representation can be broadly divided into statistical and structural methods [6]. In the former, the document is represented by a feature vector, and in the latter, a data structure (e.g. graphs or trees) are used to describe objects and their relationships in the document. The classical retrieval systems are often limited to work with a statistical representation due to the need of computing distances between documents (feature vectors) or finding a representative of a cluster of documents. However, when a numerical feature vector is used to represent the document, all structural information is discarded although the structural representation is more powerful in terms of its representational abilities [6]. In the last decades, many structural approaches have been proposed. These approaches deal, especially, with graph-based representations. Nevertheless, dealing with graphs suffers, on the one hand from the high complexity of the graph matching problem which is a problem of computing distances between graphs, and on the other hand from the robustness to structural noise which is a problem related to the capability to cope with structural variations and differences in the size of the graph. In order to overcome this problem, several approximate graph matching methods have been proposed [13, 18, 22, 24]. In this paper, our attention is focused on the comparison of different graph similarity measures in the context of document retrieval.

Graph similarity measures use different techniques to minimize the complexity and to optimize the robustness to structural noise. Robles-Kelly and al. [24] propose a spectral seriation approach to reduce the graph matching to a string edit distance in a probabilistic framework. Jouili and al. [13] simplify the problem to a bipartite graph matching problem by making use of node signatures. Lopresti and al. [18] use a probe

This work is partially supported by the French National Research Agency project NAVIDOMASS referenced under ANR-06-MCDA-012 and Lorraine region.

technique to reduce the graph matching to distance between vectors. Papadopoulos and al. [22] introduce an histogram-based technique.

In this paper, we present an evaluation of these four graph distance measures on four different document data sets. We use the well-known GREC [23, 9] data base which consists of graphs representing symbols from architectural and electronic drawings. Here the ending points (ie corners, intersections and circles) are represented by nodes which are connected by undirected edges and labeled as lines or arcs. We have also performed a retrieval evaluation on an ornamental letters data set which contains lettrine (graphical object) extracted from digitized ancient document ¹. Since one lettrine contains a lot of information (i.e. texture, decorated background, letters), the graphs are extracted from a region-based segmentation [11] of the lettrine with a user-based parameterization technique. The nodes of the graph are represented by the regions and the edges describe their adjacency relationships. We have also evaluated the graph similarity measures on a shape data set [26] in which the graph is extracted by making use of a skeletonizing algorithm and a delaunay triangulation of detected endpoints. Finally, the graph similarity measures are evaluated on a set of trademark-logos in which the graph is extracted by making use of an interest point detector [12] and the delaunay triangulation.

The performance evaluation is performed using the ROC curves. Through this evaluation, we will examine the robustness of each graph similarity distance and this will allow us to investigate the applicability of each measure to the problem of retrieval for different kinds of documents.

2 Graph-based representations

Region-Based approaches have been one of the most important research issues in content-based image retrieval. Representing images at the region level captures not only the local variations of regions but also their spatial organizations. Graph-based representations are widely used in region-based segmentation. To incorporate both region attributes and adjacent relationship an image is usually represented as an attributed graph. Classical image representations such as colors histograms, texture descriptors, or shape descriptions ignore the regions localization in the document.

Graph-based representations are used in many applications, for instance, to represent circuit diagrams [4], for shape recognition [8], image matching [17, 2], or old document analysis [14]. Other works on graph-based representation [1, 3, 21, 20, 10], use different methods to incorporate features of the document image. The methods vary according to the characteristics of the data and the aims of the representation (i.e. matching or retrieval). Bunke [4] illustrates an example of converting a circuit diagram to a graph by representing the lines in the circuit diagram; each graph node represents a line endpoint, corner or intersection point, and node attributes record the image coordinates (x,y) of this feature. In [14] the authors manipulate initial letters from old documents. They proceed by segmenting the initial letter into different information layers to obtain "Information layers of homogeneous zones". Then, each homogeneous zone of the initial letter is converted to a node of graph with two attributes: size and shape descriptions, and each edge contains two attributes: angle and distance. Baeza-Yates and al. [2] also represent images as attributed graphs and adopt the graph edit distance to calculate the image distance. In another way use graphs in image analysis, Pan and al. [21] introduce a graph-based automatic image annotation. The authors propose a graph-based method to assign automatically keywords to an image. The main idea of this work is to represent all the images, as well as their attributes (caption words and regions) as nodes and link them according to their known association into a graph. For the task of image annotation, they use a "3-layer" graph, with one layer for image nodes, one layer for annotation term nodes, and one layer for the image regions.

In this section, we have seen that graphs can be widely used as a data structure-model in the pattern recognition domain. Moreover, most of the previous graph-based representations aim to measure some similarity

¹Provided by the CESR - University of Tours on the context of the ANR Navidomass project <http://l3iexp.univ-lr.fr/navidomass/>

between objects in a further recognition or retrieval task. This fact leads to the development of several similarity measures for graphs.

3 Graph matching measures

An important step in structural pattern recognition is the representation of documents by a graph data structure. This structural representation should provide a description of the characteristics of the images efficient for the task under consideration (e.g. retrieval). The retrieval problem can then be addressed in the corresponding graph space without addressing to the original images. The process of comparing graphs is generally referred as graph matching. Generally, given two graphs $G_1=(V_1,E_1)$ and $G_2=(V_2,E_2)$, the graph matching methods are divided into two broad categories: the first contains exact matching methods called graph isomorphism that requires to find a one-to-one mapping $f:V_1 \mapsto V_2$ such that $(u,v) \in E_1$ if $(f(u),f(v)) \in E_2$ with $|V_1|=|V_2|$. The second category contains inexact matching methods, where a strict correspondence among the nodes or the edges of the two graphs can not be found. Therefore, in these cases no isomorphism can be expected between both graphs, and the graph matching does not consist in finding the *exact* matching but the *best* matching between them. To perform such a structural matching, various formalisms have been proposed, using error-tolerant methods based on continuous optimization [19], quadratic programming, and spectral decomposition of graph matrices [24]. Other methods try to characterize the properties of graphs using a vector-based representation in order to profit from the existing vector measures[18, 13, 22]. Most of the inexact graph matching measures are based on some sort of edit operations. The basic idea is to define the similarity of graphs based on the effort needed to make the graphs identical. This is an extension of the well known string edit distance [16] to the graph edit distance (GED) [25]. For a review of graph similarity measures we refer the readers to [7, 5].

The matching methods selected for our evaluation belong to different formalisms. The spectral technique proposed by the Robles-Kelly's method [24] has proven to obtain good performance results. The graph matching based on node signature [13] uses a local decomposition of graphs and the assignment problem to carry out an optimum node-to-node correspondence. Papadopoulos and al. [22] provide a histogram-based representation for graphs to compute the edit distance between graphs as a sequence of three different primitive operations. Finally, using the new concept of probe, Lopresti [18] introduces the graph probing which is characterized by its rapidity.

3.1 Graph edit distance from spectral seriation

Robles-Kelly and al. [24] use a spectral method to represent graphs by strings, and then the similarity of graphs is measured according to the edit distance of strings in a probabilistic framework. The graph edit distance is the cost of the shortest edit path in an edit lattice for transforming the data graph into the model. The rows and columns of edit lattice are indexed by two strings $Y=\{y_1,y_2,\dots,y_{|V_D|}\}$ for data graph $G_D=(V_D,E_D)$ and $X=\{x_1,x_2,\dots,x_{|V_M|}\}$ for the model graph $G_M=(V_M,E_M)$, with null symbol ε , and V_D and E_D being the point set and the edge set of the data graph. The problem of computing the edit distance is posed as that of finding the least expensive path $\Gamma^* = \langle \gamma_1, \gamma_2, \dots, \gamma_k, \gamma_L \rangle$ from (y_1, x_1) to $(y_{|V_D|}, x_{|V_M|})$ through the edit lattice based on the Levenshtein distance. Each state $\gamma_k \in (V_D \cup \varepsilon) \times (V_M \cup \varepsilon)$ of the edit path is a Cartesian pair. Then cost functions are defined for elementary matches, according to the cost edit path Γ^* (i.e., graph edit distance) computed using the following equation:

$$d(X, Y) = C(\Gamma^*) = \sum_{\gamma_k \in \Gamma} \eta(\gamma_k \rightarrow \gamma_{k+1}) \quad (1)$$

where $\eta(\gamma_k \rightarrow \gamma_{k+1}) = -(\ln P(\gamma_k | \phi_X^*(x_i), \phi_Y^*(y_j)) + \ln P(\gamma_{k+1} | \phi_X^*(x_{i+1}), \phi_Y^*(y_{j+1})) + \ln R_{k,k+1})$, and the edge compatibility coefficient $R_{k,k+1}$ is

$$R_{k,k+1} = \frac{P(\gamma_k|\gamma_{k+1})}{P(\gamma_k)P(\gamma_{k+1})} = \begin{cases} \rho_M \rho_D & \text{if } \gamma_k \rightarrow \gamma_{k+1} \text{ is a diagonal transition on the edit lattice} \\ \rho_M & \text{if } \gamma_k \rightarrow \gamma_{k+1} \text{ is a vertical transition on the edit lattice} \\ \rho_D & \text{if } \gamma_k \rightarrow \gamma_{k+1} \text{ is a horizontal transition on the edit lattice} \\ 1 & \text{if } y_j = \varepsilon \text{ or } x_i = \varepsilon \text{ and } y_{j+1} = \varepsilon \text{ or } x_{i+1} = \varepsilon \end{cases}$$

$$P(\gamma_k|\phi_X^*(x_i), \phi_Y^*(y_j)) = \begin{cases} \frac{1}{\sqrt{2\pi\sigma}} \exp\{-\frac{1}{2\sigma^2}(\phi_X^*(x_i) - \phi_Y^*(y_j))^2\} & \text{if } y_j \neq \varepsilon \text{ and } x_i \neq \varepsilon \\ \alpha & \text{if } y_j = \varepsilon \text{ and } x_i \neq \varepsilon \end{cases}$$

where ρ_M and ρ_D are respectively the edge densities of the graphs G_M and G_D ($\rho_M = \frac{|V_M|^2}{E_M}$) and ϕ_X^* and ϕ_Y^* are, respectively, the leading eigenvectors of the adjacency matrices for the graph G_M and G_D . In the remainder, we denote this graph matching technique by GEDSS.

3.2 Graph matching based node signatures

Jouili and al. [13] propose a new algorithm for matching and computing the distance between weighted graphs. They introduce a new *vector-based node signature*. Each node is associated with a vector $\gamma(x) = \{d(x), w_0, w_1, w_2 \dots\}$, where x is a node, $d(x)$ gives the degree of x , and w_i are the weights of the incident edges to x . Using these node signatures a cost matrix C is constructed. The cost matrix $C_{g_i, g_j}(i, j) = L_1(\gamma(i), \gamma(j))$ describes the matching costs between nodes in two graphs. It is a (n, m) matrix where n and m are the sizes of the two graphs. An element (i, j) in this matrix gives the Manhattan distance between the i th node signature in the first graph and the j th node signature in the second graph. To find the optimum matching, the problem is considered as an instance of the assignment problem, which can be solved by the Hungarian method [15]. They define the distance between two graphs g_i and g_j as follows:

$$D(g_i, g_j) = \frac{\hat{M}}{|M|} + ||g_i| - |g_j|| \quad (2)$$

where

- $|M|$: the size of the matching function M which is the number of matching operations. In any case, when two graphs are matched the number of the matching operations is the size of the smaller one.
- $\hat{M} = \sum L_1(\gamma(x), \gamma(M(x)))$: the matching cost which is the sum of the matching operation costs, for two graphs matched by M .

In the remainder, we denote this graph matching technique by GMNS.

3.3 Graph probing approach

Lopresti and al. [18] introduce the paradigm of graph probing. This technique consist on using a probe into the graphs to determine some particular information. The measure of similarity between two graphs is an L1 norm distance of the two corresponding vectors. For the construction of vectors, Lopresti present three classes of construction each one led by a question, Class 0: "How many vertices with degree n are present in graph $G = (V, E)$?", Class 1: "How many vertices with in-degree m and out-degree n are present in G ?", Class 2: "How many vertices labeled att are present in G ?". The use of such class depends on the type of graph. Let $G = (V, E)$ be an undirected graph, the vector associated to G is: $PR(G) \equiv (n_0, n_1, n_2, \dots)$ where $n_i = |\{v \text{ in } V \mid \deg(v) = i\}|$. So the distance between two graphs is $L_1(PR1, PR2)$. In the remainder, we denote the graph probing technique by GP.

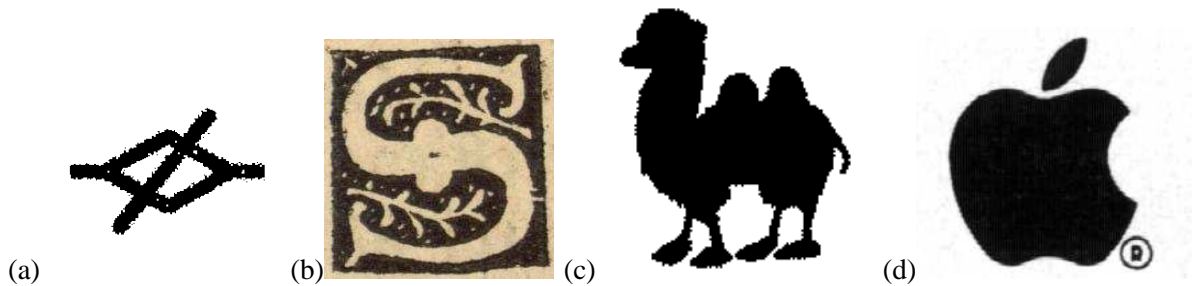


Figure 1: Samples from: (a) GREC database, (b) Ornamental letters database, (c) Shape database, (d) Logo database

3.4 Graph histogram approach

Papadopoulos and al. [22] present a similarity measure for graphs, which is based on the concept of edit operations. They propose three different primitive operations, which are vertex insertion, vertex deletion and vertex update. While vertex insertions or deletions have a trivial meaning, the update operation is needed to insert or delete edges incident to a vertex. Additionally they introduce the degree sequence of a graph, i.e. the non-increasing sequence of the degrees of vertices in a graph. The similarity distance between two graphs is defined as the minimum number of primitive operations which are required so that the two graphs have the same degree sequence. To calculate the similarity measure, the sorted graph histogram is introduced, which is a histogram of the degrees of the vertices in a graph. Papadopoulos and al. show also that the L_1 -distance between two sorted graph histograms defines their similarity distance. Additionally it is proven that the similarity distance satisfies the metric properties. In some cases, the sorted degree histograms of the graphs in a database are of different dimensionality if not all graphs are of the same order. To allow the use of index structures for vector spaces, the authors introduce a histogram folding technique to achieve a constant dimensionality of the histograms for all graphs. In the remainder, we denote this method by GH.

4 Experimental results

4.1 Data sets

The graph retrieval tasks considered in this paper include the retrieval of line drawings (symbols), ancient documents (ornamental letters), LEMS's shape database and the set of trademark-logos. Figure 1 represents samples for each database used.

- **GREC database:** The GREC database [23, 9] (see figure 1(a)) consists of graphs representing symbols from architectural and electronic drawings. Here, the ending points (ie corners, intersections and circles) are represented by nodes which are connected by undirected edges and labeled as lines or arcs. The graph database used in our experiments has of 528 graphs, 24 classes and 22 graphs per class.
- **Ornamental letters database:** The ornamental letters database (see figure 1(b)) contains lettrine (graphical object) extracted from digitized ancient document ². Since one lettrine contains a lot of information (i.e. texture, decorated background, letters), the graphs are extracted from a region-based segmentation [11] of the lettrine with a user-based parameterization technique. The nodes of the graph are represented by the regions and the edges describe their adjacency relationships. The graph database used in our experiments consists 280 graphs, 4 classes and 70 graphs per class.

²Provided by the CESR - University of Tours on the context of the ANR Navidomass project <http://l3iexp.univ-lr.fr/navidomass/>

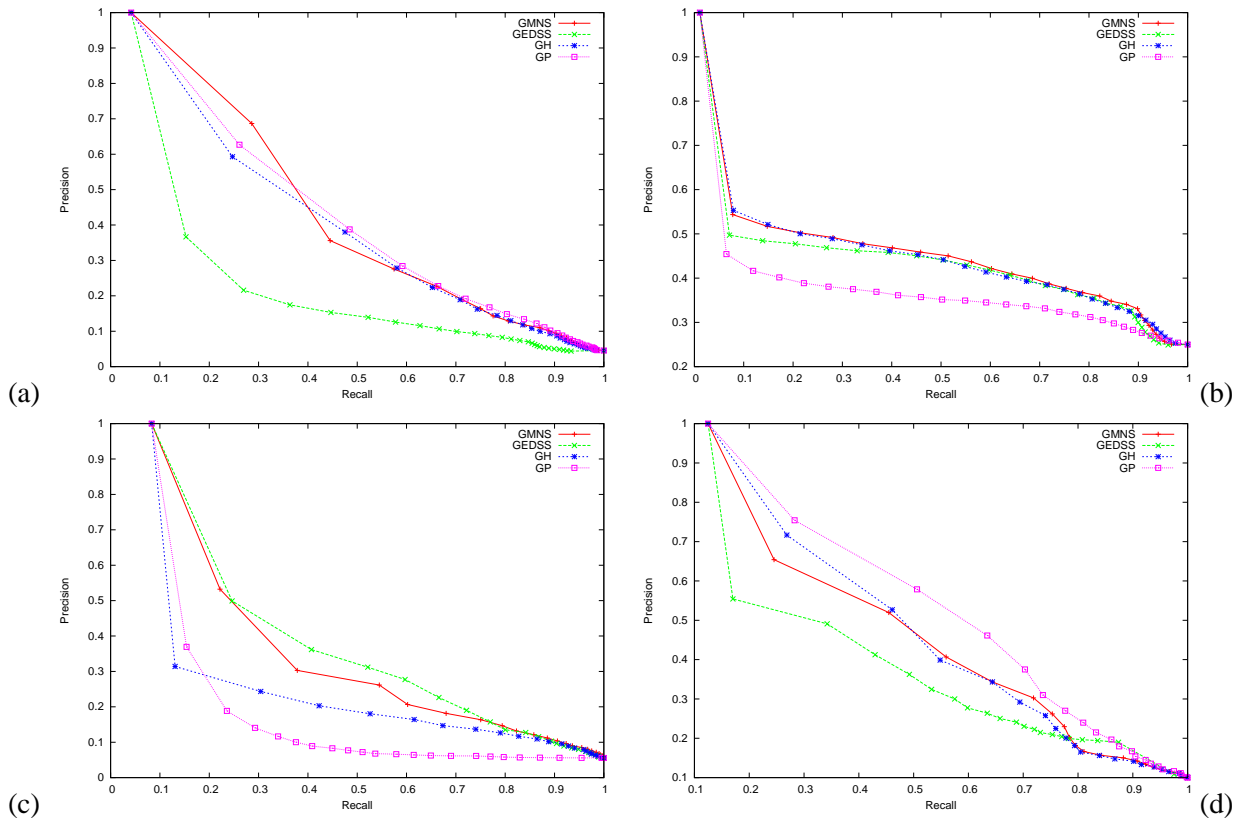


Figure 2: Precision-Recall curves on: (a) GREC database, (b) Ornamental letters database, (c) Shape database, (d) Logo database

- Shape database:** We use the shapes provided by the LEMS laboratory of the Brown University [26] (see figure 1(c)). the graphs are extracted from the shapes by skeletonizing and applying a polygonal approximation to the skeleton to obtain straight line segments. For each line segment, we locate endpoints and the graphs are based on the Delaunay triangulations of these endpoints. The graph database used in our experiments has 216 graphs, 18 classes and 12 graphs per class.
- Logo database:** This database (see figure 1(d)) consists of graphs representing binary images of trademark-logos. Here, graphs are extracted by the delaunay triangulations of the detecting points of interest by the Harris algorithm [12]. The graph database used in our experiments consists of 80 graphs, with 10 classes and 8 graphs per class.

4.2 Results analysis

The results of our experiments for these four databases with the four graph matching measures are presented in figure 2.

From the precision-recall curves, we can remark that the performance of the graph matching methods depend on the databases. For the GREC database, the matching measures (GP, GH and GMNS) that use simple structural modification perform similarly and better than the GEDSS method which use a string representation for graphs. We realize that for graphs with low edge and node densities (as the case of the GREC database) the string-based representation is not discriminant. In addition, the GMNS method provides a performance peak for low recall values, and it joins the performance of the GP and GH methods for high recall values. The discrimination of the node signatures provides a good robustness for this kind of database.

From the results provided on the Lettrine database, we see that all the distance measures provide similar results with a little less performance for the GP technique. This may be explained by the fact that the different methods produce a quite similar response to the structural errors between the graphs used to represent the ornamental letters. In the other way, one can conclude that this kind of graph representation (region adjacency graph) of the ornamental letter is more or less robust to different graph matching methods.

In the case of the shape database, the performance of the graph probing fails clearly in comparison with other distance measures. It seems that the probe of the node degree is not a good discriminating feature for this database which presents important structural errors between graphs in different classes. Further, the GEDSS method which has shown previously good results for similar databases (see [24]), provides the better retrieval results.

For the logo database, all the distance measures provide similar behaviors. Here, the graph probing keeps the leader position among the other distances. In addition, the provided results of all the distance measures are particularly better in comparison with the other databases. This may be due to the suitable graph representation used for this database. We can think that the graph representation approaches used for other databases is not necessary the most suitable. In addition, different distance measures provide quite similar results for a given graph representation as the case of the Ornamental letters database. From all these results, we can remark that the GP and GEDSS methods are more sensitive to the representation we put in the graph.

5 Conclusion and perspectives

In this paper we have compared the performance of four graph matching methods for graph retrieval with different kind of document databases. The receiver-operating curve (ROC) is used to measure retrieval performances. The ROC curve is formed by Precision rate against Recall rate. Our experimental results show that the performance of each graph distance measure depends on the databases. That is to say, a given graph distance can provide a good performance for one database and poor performance for an other database. Moreover, for a good graph representation we can remark that the performances of different graph matching methods are quite similar. Approaches are also more and less robust to the variability of the representation. In future works we want to study the behavior of these methods against the representation we put in the graph and the type of database.

References

- [1] R. Ambauen, S. Fisher, and B. H. Graph edit distance with node splitting and merging, and its application to diatom identification. *IAPR-TC15 Workshop on GbRPR, LNCS 2726*, pages 95–106, 2003.
- [2] R. Baeza-Yates and G. Valiente. An image similarity measure based on graph matching. *Proc. International Symposium on String Processing Information Retrieval*, pages 28–38, 2000.
- [3] H. G. Barrow and R. J. Popplestone. Relational descriptions in picture processing. *In Machine Intelligence*, 4:377–396, 1971.
- [4] H. Bunke. Attributed of programmed graph grammars and their application to schematic diagram interpretation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 4(6):574–582, Novembre 1982.
- [5] H. Bunke. Recent developments in graph matching. *Pattern Recognition, 2000. Proceedings. 15th International Conference on*, pages 117–124 vol.2, 2000.
- [6] H. Bunke, S. Günter, and X. Jiang. Towards bridging the gap between statistical and structural pattern recognition: Two new concepts in graph matching. In S. Singh, N. A. Murshed, and W. G. Kropatsch, editors, *ICAPR*, volume 2013 of *Lecture Notes in Computer Science*, pages 1–11. Springer, 2001.

- [7] D. Conte, P. Foggia, C. Sansone, and M. Vento. Thirty years of graph matching in pattern recognition. *International Journal of Pattern Recognition and Artificial Intelligence*, 18(3):265–298, 2004.
- [8] C. Di Ruberto, G. Rodriguez, and L. Casta. Recognition of shapes by morphological attributed relational graphs. *AIIA*, 2002.
- [9] P. Dosch and E. Valveny. Report on the second symbol recognition contest. In W. Liu and J. Lladós, editors, *GREC*, volume 3926 of *Lecture Notes in Computer Science*, pages 381–397. Springer, 2005.
- [10] J. Fauqueur and N. Boujemaa. Region-based image retrieval: fast coarse segmentation and fine color description. *J. Vis. Lang. Comput.*, 15(1):69–95, 2004.
- [11] P. F. Felzenszwalb and D. P. Huttenlocher. Efficient graph-based image segmentation. *International Journal of Computer Vision*, 59(2):167–181, 2004.
- [12] C. Harris and M. Stephens. A combined corner and edge detection. In *Proceedings of The Fourth Alvey Vision Conference*, pages 147–151, 1988.
- [13] S. Jouili and S. Tabbone. Graph matching using node signatures. In *IAPR-TC15 Workshop on GbrPR, Italy, LNCS 5534, Springer*, pages 154–163, 2009.
- [14] A. Karray, J.-M. Ogier, S. Kanoun, and M. A. Alimi. An ancient graphic documents indexing method based on spatial similarity. In W. Liu, J. Lladós, and J.-M. Ogier, editors, *GREC*, volume 5046 of *Lecture Notes in Computer Science*, pages 126–134. Springer, 2007.
- [15] H. W. Kuhn. The hungarian method for the assignment problem. *Naval Research Logistics Quarterly*, 2:83–97, 1955.
- [16] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966.
- [17] C.-Y. Li and C.-T. Hsu. Region correspondence for image retrieval using graph-theoretic approach and maximum likelihood estimation. *International Conference on Image Processing*, pages 421–424, 2004.
- [18] D. P. Lopresti and G. T. Wilfong. A fast technique for comparing graph representations with applications to performance evaluation. *IJDAR*, 6(4):219–229, 2003.
- [19] M. Neuhaus. Edit distance-based kernel functions for structural pattern classification. *Pattern Recognition*, 2006.
- [20] I. Ounis and M. Pasca. Modeling, indexing and retrieving images using conceptual graphs. In G. Quirchmayr, E. Schweighofer, and T. J. M. Bench-Capon, editors, *DEXA*, volume 1460 of *Lecture Notes in Computer Science*, pages 226–239. Springer, 1998.
- [21] J. Pan, H. Yang, C. Faloutsos, and P. Duygulu. Gcap: Graph-based automatic image captioning. In *Proceedings of the 4th International Workshop on Multimedia Data and Document Engineering*, 2004.
- [22] A. N. Papadopoulos and Y. Manolopoulos. Structure-based similarity search with graph histograms. *Proceedings of International Workshop on Similarity Search (DEXA IWOSS'99)*, pages 174–178, September 1999.
- [23] K. Riesen and H. Bunke. Iam graph database repository for graph based pattern recognition and machine learning. In N. da Vitoria Lobo, T. Kasparis, F. Roli, J. T.-Y. Kwok, M. Georgiopoulos, G. C. Anagnostopoulos, and M. Loog, editors, *SSPR/SPR*, volume 5342 of *Lecture Notes in Computer Science*, pages 287–297. Springer, 2008.

- [24] A. Robles-Kelly and E. R. Hancock. Graph edit distance from spectral seriation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 27(3):365–378, 2005.
- [25] A. Sanfeliu and K. Fu. A distance measure between attributed relational graphs for pattern recognition. *IEEE Trans. on Systems, Man and Cybernetics*, 13(3):353–362, May 1983.
- [26] D. Sharvit, J. Chan, H. Tek, and B. B. Kimia. Symmetry-based indexing of image databases. *J. Visual Communication and Image Representation*, 9:366–380, 1998.