



# AnyTime and Distributed Approaches for Graph Matching

**Presented by**: Zeina Abu-Aisheh

Wednesday, May 25 2016

Supervised by:

Prof. Jean-Yves Ramel, Dr. Romain Raveaux and Prof. Patrick Martineau

### **External Jury Members:**

Prof. Christine Solnon Prof. Francesc Serratosa

Prof. Luc Brun Dr. Kasper Riesen

# Surrounded by Graphs



 $\mathbf{*}$ Introduction

Strength and Weakness of Graph Matching Methods

Contribution1: Optimization Techniques

**Contribution2:** Metrics and Datasets for Performance Evaluation

\* Experiments

 $\mathbf{Conclusion}$ 



## Outline

#### **\*Introduction**

Graph Matching Problems

Strength and Weakness of Graph Matching Methods

**Contribution 1: Optimization Techniques** 

**\***Contribution 2: Metrics and Datasets for Performance Evaluation

\* Experiments

 $\diamond$  Conclusion

## Exact Graph Matching

### Graph Isomorphism



## Subgraph Isomorphism



- Strict correspondances
  - ✓ Edge-preserving
- Too strict not used in
- Pattern Recognition

## Error-Tolerant Graph Matching

- Graphs of different topology/structures and attributes can be matched
- Attributes penality
- Structure penality

• NP-complete

• Graph Isomorphism "open problem"

### Graph Edit Distance (GED)



Let  $G_1 = (V_1, E_1, L_{V_1}, L_{E_1}, \mu_1, \zeta_1)$  and  $G_2 = (V_2, E_2, L_{V_2}, L_{E_2}, \mu_2, \zeta_2)$  be two graphs, the graph edit distance between  $G_1$  and  $G_2$  is defined as:

$$d_{\lambda_{min}}(G_1, G_2) = \min_{\lambda \in \gamma(G_1, G_2)} \sum_{ed_i \in \lambda} c(ed_i)$$



## Edge Mappings based on their Adjacent Vertices



### Focus on Graph Edit Distance (GED)

#### • Cost Functions:

✓ GED is generic and compatible with graphs labeled with numeric and/or symbolic attributes

 $\checkmark$  GED can be a metric (under constraints)

 $\checkmark$  GED can be turned into the Maximum Common Subgraph problem

# Outline

Introduction

### Strength and Weakness of Graph Matching Methods

- Graph Edit Distance Methods
- Graph Matching Repositories and their performance evaluation metrics
- Limitations of existing methods

**\***Strength and Weakness of Graph Matching Methods

**\***Optimization Techniques

**\*** Metrics and Datasets for Performance Evaluation

\* Experiments

#### Conclusion







## Graph Edit Distance Approaches





• Tree-Based Search Approaches











- Tree-Based Search Approaches
  - ✓ Mathematical Formulation [Justice and Hero 2006]



It can only match unattributed edges



## Graph Edit Distance Approaches





- Search Space Truncation
- Problem Reformulation



- Search Space Truncation
  - ✓ Beam Search [Riesen 2006]
- Problem Reformulation

At each iteration, we keep the *x* most promising solutions
✓ Suboptimal solution



- Search Space Truncation
- Problem Reformulation
  - ✓ Bipartite Graph Matching: [Riesen 2009]





- Search Space Truncation
- Problem Reformulation
  - ✓ Fast Bipartite Graph Matching; [Serratosa 2014]

	$c_{1,1}$			$c_{1, V_2 }$	$c_{1,\epsilon}$	$\infty$		$\infty$
$C_{ve} =$					$\infty$	$c_{2,\epsilon}$		$\infty$
	$c_{ V_1 ,1}$			$c_{ V_1 , V_2 }$	$\infty$		$\infty$	$c_{ V_1 ,\epsilon}$
	$c_{\epsilon,1}$	$\infty$		$\infty$	0			0
	$\infty$	$c_{\epsilon,2}$		$\infty$				
	$\infty$		$\infty$	$C_{c}  V_{c} $	0			0

$$C_{ve} = \begin{vmatrix} c_{1,1} - (c_{1,\epsilon} + c_{\epsilon,1}) & \dots & \dots & c_{1,|V_2|} - (c_{1,\epsilon} + c_{\epsilon,|V_2|}) \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ c_{|V_1|,1} - (c_{|V_1|,\epsilon} + c_{\epsilon,1}) & \dots & \dots & c_{|V_1|,|V_2|} - (c_{|V_1|,\epsilon} + c_{\epsilon,|V_2|}) \end{vmatrix}$$



- Search Space Truncation
- Problem Reformulation
  - ✓ Bipartite Graph Matching: [Riesen 2009]
    - Local rather than global structures
  - $\checkmark$  Improvements by Local Search





- Search Space Truncation
- Problem Reformulation
  - ✓ Bipartite Graph Matching: [Riesen 2009]
    - Local rather than global structures
  - ✤ Improvements via Local Search Strategies
    - ➢ Genetic algorithm: [Riesen 2014]
    - Beam Search algorithms:
      - □ BP-Beam: [Riesen 2014] in Artifical Neural Network in PR
      - □ IBP-Beam: [Ferrer 2015] in GBR2015
      - □ SBP-Beam: [Ferrer 2015] in PRL



## Limitations of Graph Edit Distance Approaches

- No intense study on the exact graph edit distance side.
  - $\checkmark$  Long run time and high memory consumption
  - $\checkmark$  Good to evaluate accuracy
- Approximate methods are fast enough, BUT:
  - $\checkmark$  No deep study on the quality of the solution as a function of time
- No algorithms work under time and memory constraints
  - $\checkmark$  Tree-based algorithms can be of great interest "anytime algorithms"

**Contribution**: Anytime GED that can find a good compromise between quality and speed



## Graph Matching Repositories and their Metrics

	Problem Type	Graph Type	Database Type	Metrics Type	Purpose
[Santo2003]	Exact GM	Non- attributed	Synthetic	Accuracy and scalability	Matching
CMU	<b>MU</b> Error-tolerant GM		Real-world	Memory consumption, accuracy and matching	Matching
IAM DB [Riesen 2008]	Error-tolerant GM	Attributed	Real-world	Accuracy and running time	Classification
[Foggia 2001]	Exact GM	Attributed	Synthetic	Accuracy and scalability	Matching
[Carletti 2014]	[Carletti 2014] Exact GM		Real-world	Scalability	Matching
[Pascal 2007]	Error-Tolerant GM	Attributed	Real-world	Accuracy and matching	Classification/ Matching
Tarragona in [Cortes 2015]	Error-Tolerant GM	Attributed	Real-World	Accuracy and matching	Matching



## Limitations of Existing Repositories

#### Lack of metrics for GM methods:

 $\checkmark$  None of them studied the scalability of methods

✓ Lack of diversity of graph datasets (dense, large graphs, etc.)

<u>**Contribution:**</u> An annotated repository dedicated to scalability with performance evaluation metrics

#### Introduction

\*Strength and Weakness of Graph Matching Methods

### **Contribution 1: Optimization Techniques**

- Toward an Anytime Graph Matching Approach
- Parallel and Distributed Approaches for Graph Edit Distance

**\***Contribution 2: Metrics and Datasets for Performance Evaluation

**\***Experiments

Conclusion



## From Exact to Anytime Algorithms





From Exact to Anytime Algorithms

• Anytime Algorithms:

✓ Search-tree based algorithms are anytime algorithms by nature





## Anytime Graph Matching





From Exact Graph Edit Distance to Anytime One

### Overcoming the memory bottleneck of A\*

✓ Depth-First Graph Edit Distance (DF)



Z. Abu-Aisheh, R. Raveaux, J-Y Ramel and P. Martineau : An Exact Graph Edit Distance Algorithm for Solving Pattern Recognition Problems. ICPRAM (1) 2015 : 271-278.



## From Exact Graph Edit Distance to Anytime One





## From Exact Graph Edit Distance to Anytime One





## From Exact Graph Edit Distance to Anytime One

#### Overcoming the memory bottleneck of A\*



35



## From Exact Graph Edit Distance to Anytime One

• From *DF* to Anytime *DF* (*ADF*)





• **Interruptability:** At each iteration, we can stop the algorithm if it exceeds *CT* 

Z. Abu-Aisheh, R. Raveaux and J-Y Ramel: Anytime Graph Matching. Special issue of Pattern Recognition Letters (revision: 2<sup>nd</sup> round).


# ADF: Theoretical Observations

- ADF memory complexity is O(|V<sub>1</sub>|.|V<sub>2</sub>|)
  - ✓ Encountering a local minimum while exploring
  - $\checkmark$  To have more diversity + scaling up
    - Parallel and Distributed versions of ADF





# **Raised Questions**

- 1. What are the subtasks that should be generated before the parallelism starts?
- 2. How to distribute the tasks?
- 3. How to keep all threads busy?
- 4. How to prune the search space as fast as possible?



# Coming back to Anytime Depth-First (ADF)

- Tasks = exploration of partial solutions is more time consuming than g(p) and h(p)
- Search tree is irregular
  - $\checkmark\,$  Load Balancing is necessary



- Prune the search tree: sharing the upper bound
- Number of CPU/machines
  - ✓ No rule (experiments)



# Parallel and Distributed Branch-and-Bound Algorithms

• None of branch-and-bound algorithms was dedicated to solving error-tolerant graph matching





# A Parallel Depth-First Graph Edit Distance (**PDFS**)





# PDFS: Decomposition and Assignment



## Ordered list (OPEN):

<b>P</b> 1	<b>P</b> 2	<b>P</b> 3	<b>P</b> 4	<b>P</b> 5		P <sub>N</sub>
------------	------------	------------	------------	------------	--	----------------





# PDFS: Load Balancing





# Parallel and Distributed Branch-and-Bound Algorithms

• None of branch-and-bound algorithms was dedicated to solving error-tolerant graph matching.



- No load balancing
  - ✓ Fault Tolerance

Z. Abu-Aisheh, R. Raveaux, J-Y Ramel and P. Martineau: A Distributed Algorithm for Graph Edit Distance. GraphSM'2016.



# A Distributed Graph Edit Distance Algorithm (**D-DF**)





# **D-DF**: Decomposition and Assignment



#### Ordered list (OPEN):







# **PDFS** and **D-DF** Theoretical Observations

## • **PDFS** and **D-DF**:

## ✓ Inputs:

- Generate N nodes before dispatching the work
- $\succ$  *h*: estimated cost
- > PDFS: Number of threads
- > *D-DF*: Number of processes

 $\checkmark$  Exploring different parts of the search tree

## • *D-DF*:

- ✓ Single job
- $\checkmark$  Workers may be left without any task to do



Ordered list (OPEN):

# Outline

### Introduction

Strength and Weakness of Graph Matching Methods

## **\*Optimization Techniques**

## **\* Metrics and Datasets for Performance Evaluation**

- GDR4GED: Graph Repository
- GDR4GED: Performance Evaluation Metrics

## Experiments

## Conclusion



# GDR4GED: Graph Repository - Databases

Database	Decomposition	Overview	Purpose
GREC	MIX, 5, 10, 15 and 20		Classification
MUTA	MIX ,10, 20, 30, and 70		Classification
Protein	MIX, 20, 30 and 40	And the second second	Classification
CMU	30		Matching

**Repository link**: https://iapr-tc15.greyc.fr/links.html#Benchmarking and data sets

Z. Abu-Aisheh, R. Raveaux and J-Y Ramel: A Graph Database Repository and Performance Evaluation Metrics for Graph Edit Distance. GbRPR 2015 : 138-147.



# GDR4GED: Graph Repository – Graph-Level Annotations

## • 7 methods (Time constraint: 300 seconds per comparison):

## $\checkmark$ Exact methods:

- ≻ *A*\* [Riesen 2007]
- ≻ JHBLP [Justice 2006]
- ➢ DF [Abu-Aisheh 2015]
- ► D-DF [Abu-Aisheh 2016]

## $\checkmark$ Approximate methods:

- ➢ BeamSearch [Neuhaus 2006]
- ▶ BP [Riesen 2009]
- ≻ Hausdorff [Fischer 2013]

## • For each graph matching pair:

- $\checkmark\,$  The name of the method that found the best solution
  - > Distance between each pair of graphs (Optimal/suboptimal)
  - Vertex-Vertex matching



# **GDR4GED:** Performance Evaluation Metrics

• Under time and memory constraints





# Metrics for each Pair of Graphs

• Given the reference R:

- ✓ Deviation
- ✓ Matching Dissimilarity

✓ Best found solution

- Reference Distance (RD) = 5
- Reference Matching (RM):

Distance (m) = 6 ✓ Deviation = 20%

٠

- Matching Dissimilarity = 75%
- Best found solution = ?







# Metrics for a Full Dataset

- Number of unfeasible solutions:
  - $\checkmark$  i.e. incomplete editpath
- Number of optimal solutions:
- Number of Time-Out and Out-Of-Memory cases
- Mean number of explored nodes
- Mean running time in milliseconds
- Running time-Deviation plot:
  - ✓ Projection on a two-dimensional space ( $\mathbb{R}^2$ )

# Outline

 $\mathbf{*}$  Introduction

\* Strength and Weakness of Graph Matching Methods

**\*** Optimization Techniques

**\*** Metrics and Datasets for Performance Evaluation

#### **\*** Experiments

- Included Methods
- Tests under Soft and Hard Time Constraints
- Quality as a Function of Time
- Density Tests
- Classification under Time and Memory Constraints

### Conclusion



# **Included Methods**

	Methods	Reference	Parameters
_	<b>A</b> *	[Riesen 2007]	_
Exact _	JHBLP	[Justice 2006]	_
	BP	[Riesen 2009]	-
Approximate	FBP	[Serratosa 2014]	_
	Н	[Fischer 2013]	_
	BS-x	[Riesen 2006]	<i>x</i> = 1, 10 and 100



# Our Exact Approaches and their Parameters

Sequential Methods, ADF and PDFS:
✓ 24-core Intel i5 processor 2.10GHz
16GB memory

#### • *D-DF*:

 ✓ 5-node cluster of machines running Hadoop MapReduce version 1.0.4
Each node contains a 4-core Intel i7
processor 3.07GHz, 8GB memory
✓ Zookeeper: Message passing tool

Methods	Parameters	
ADF	<i>h(p)</i> = heursitic <i>(BP)</i>	
PDFS	h(p)= heursitic(BP) # threads = 64 N = 100	
D-DF	Lb(p)= heursitic(BP) # machines = 5 # workers per machine = 4 N = 250	

# Objectives of the Experiments

## 1. Tests under soft and hard time constraints.

- ✓ Soft Constraint: 300 seconds per comparison
- ✓ Hard Constraint:

$$C_T = \max_{m,i,j} \{ time_m(G_i, G_j) \}$$

m is BS, BP or FBP

- 2. Trade-off between quality and time
- 3. Classification under time constraints



## Soft Time Constraints: MUTA – 300 seconds



58



# Soft Time Constraints: MUTA – 300 seconds





# Soft Time Constraints: CMU – 300 seconds

• Out of memory

✓ *BS-100* and  $A^*$ 

- *BS-10* is the worst
- *BS-1* is the best ✓ 10.6%

• *D-DF* (0.78%)





# Hard Time Constraints

- *D-DF* cannot be tested under hard time constraints
- Since *C*<sup>*T*</sup> is quite small:
  - $\checkmark$  The number of threads in *PDFS* is set to 3
  - $\checkmark$  Removing the lower bound from all of *BS*, *A\**, *ADF* and *PDFS*
  - $\checkmark N$  is limited to the first level of the search tree





# Hard Time Constraints: MUTA – 500 milliseconds

- JHBLP
  - ✓ unfeasible solutions for Graphs > 30 vertices
- **PDFS:** 18%
- **BS-100** outputted unfeasible solutions on MUTA-50, 60, 70 and MIX
- **BP**: 40.37%





# Hard Time Constraints: CMU – 500 milliseconds

• *BS-10* and *BS-100* outputted unfeasible solutions





# Deviation as a Function of Time – 40 milliseconds







# Deviation as a Function of Time – 500 milliseconds







# **Density Scalability**

Why ADF was faster on CMU?

- FBP and BP are faster when having low density graphs
- ADF is faster when having high density graphs
- Studying the behavior of GED methods when increasing density
  - ✓ Setup Time
  - ✓ Four subsets (0.1, 0.2, 0.4 and 0.6)
  - $\checkmark$  10 graphs of 100 vertices per subset, each subset's size is 100 vertices



## ADF versus FBP and BP



- 1-NN Classifier
  - ✓ GREC: train set (286 graphs), test set (528 graphs) -- 400 milliseconds
  - ✓ Protein: train set (200 graphs), test set (200 graphs) -- 400 milliseconds
  - ✓ Mutagenicity: train set (1500 graphs), test set (2337 graphs) -- 500 milliseconds



- 1-NN Classifier
  - ✓ GREC: train set (286 graphs), test set (528 graphs) -- 400 milliseconds
  - ✓ **Protein:** train set (200 graphs), test set (200 graphs) -- 400 milliseconds
  - ✓ Mutagenicity: train set (1500 graphs), test set (2337 graphs) -- 500 milliseconds

Methods	<b>Classification Rate</b>	Response Time (ms)	
ADF-UB-LB	0.985	140525.48	
PDFS-UB-LB	0.985	99850.79	
A*-LB	0.530	222045.94	
BS-1	0.985	69236.34	
BP	0.985	62294.60	
FBP	0.985	27922.65	
Н	0.962	63563.74	



## • 1-NN Classifier

- ✓ **GREC:** train set (286 graphs), test set (528 graphs) -- 400 milliseconds
- ✓ **Protein:** train set (200 graphs), test set (200 graphs) -- 400 milliseconds
- ✓ Mutagenicity: train set (1500 graphs), test set (2337 graphs) -- 500 milliseconds

Methods	<b>Classification Rate</b>	Response Time (ms)	
ADF-UB	0.52	124361.61	
PDFS-UB	0.52	80038.33	
A*-LB	0.29	1065106.80	
BS-100	0.26	141265.41	
BP	0.52	59041.84	
FBP	0.385	39425.69	
Н	0.43	71990.62	



## • 1-NN Classifier

- ✓ **GREC:** train set (286 graphs), test set (528 graphs) -- 400 milliseconds
- ✓ **Protein:** train set (200 graphs), test set (200 graphs) -- 400 milliseconds
- ✓ Mutagenicity: train set (1500 graphs), test set (2337 graphs) -- 500 milliseconds

Methods	<b>Classification Rate</b>	Response Time (ms)	
ADF	0.70089	1139134.29	
PDFS	0.70089	760861.518	
A*	0.4574	856793.020	
BS-1	0.55840	1015688.00	
BP	0.70089	528546.64	
FBP	0.70089	376135.51	
Н	0.58964	525610.25	

# Outline

 $\mathbf{*}$ Introduction

Strength and Weakness of Graph Matching Methods

 $\diamond$  Optimization Techniques

**\*** Metrics and Datasets for Performance Evaluation

\* Experiments

Conclusion and Future Work


# Conclusion

- Compare graphs under time constraints
- Study deeply matching quality (graph-level)
- Focusing on Graph Edit Distance
- Adding a new methods family called «Anytime»
   ✓ Proposing a Depth-First GED method
   ✓ Converting it to Anytime (ADF)
- Extensions to *ADF*, to be able to match larger graphs with better quality:
  - ✓ Parallel Graph Edit Distance (*PDFS*):
    - Load balancing
  - ✓ Distributed Graph Edit Disrance (*D*-*DF*):
    - ➤ Master-Slave on top of Hadoop



# Conclusion: (Repository + Evaluation Metrics)

## • Public repository (GDR4GED):

 $\checkmark$  Scalability: database decomposition

✓ https://iapr-tc15.greyc.fr/links.html#Benchmarking and data sets

### Additional low level annotation:

 $\checkmark$  Best distance and vertex-vertex matchings

 $\checkmark$  7 methods (4 exact and 3 approximate)

✓ Four databases (GREC, Mutagenicity, Protein and CMU)

## • Metrics:

- $\checkmark$  Under soft and hard time constraints
- $\checkmark$  Mertics for each pair of graphs
- $\checkmark$  Metrics for a full dataset

# Conclusion: (Experiments)

Soft Constraints: *PDFS* and *D-DF* had the minimum deviation
 ✓ Except on MUTA

#### • Hard Constraints:

✓ *PDFS* is slow, however, it was always the most precise one

#### • Deviation as a Function of Time:

✓ *FBP* and *BP* were faster on less dense graphs
✓ *ADF* was faster on more dense graphs

#### Classification Tests:

✓ *PDFS* was the best exact GED method



#### • ADF:

 $\checkmark$  Better lower bounds

✓ Learn:

Lower bounds

 $\succ$  Sort vertices of G<sub>1</sub>

 $\succ$  Stop the algorithm when the quality is sufficient

#### • PDFS:

✓ Extending *PDFS* to multi-machines (MPI)

#### • *D-DF*:

 $\checkmark$  Extending *D*-*DF* to multi-jobs algorithm

# Perspectives

## • GDR4GED:

- $\checkmark$  Integrating better answers in the database (if found)
- $\checkmark$  Expanding this repository by integrating other publicly (available databases), dense graphs, etc.

## • A general perspective:

 $\checkmark$  Knowing when to use one algorithm or another

# Thanks for your attention ...

