



An Indexing Structure for Automatic Schema Matching

LIRMM

SMDB 2007, Istanbul

Fabien Duchateau, Zohra Bellahsène, Mathieu Roche

LIRMM, Université Montpellier 2

France

Mark Roantree

ISG, Dublin City University

Ireland



Outline

LIRMM

- Motivations

- BtreeMatch: an index structure for accelerating schema matching
 - Semantic Aspect
 - Performance Aspect
 - Some Performance Results

- Conclusion and Future Work



Motivations

- Finding semantic correspondences between 2 schemas still a challenging issue
- Available semi automatic matchers focus on the quality aspect of matching
- More and more large schemas, especially on the Web



Our Approach

An Index Structure for Automatic Schema Matching

- Semantic aspect
 - terminological (Levenshtein and 3grams)
 - structural (context based using cosine measure)
- Performance aspect
 - indexing structure (B-tree)



BtreeMatch: Semantic Aspect (1/4)

LIRMM

- Context of node n
 - represents the most *important* neighbour nodes of n
 - each of them is assigned a weight depending on the relationship with n

$$\omega_1(n_c, n_i) = 1 + \frac{K}{\Delta d + |\text{lev}(n_c) - \text{lev}(n_a)| + |\text{lev}(n_i) - \text{lev}(n_a)|}$$

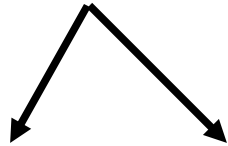
- We define StringMatching as the average between
 - Levenhstein distance
 - n-grams

BtreeMatch: Semantic Aspect (2/4)



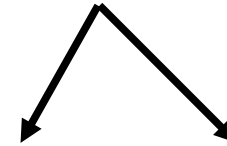
LIRMM

University



GradCourses Professor

Faculty



Courses FullProfessor

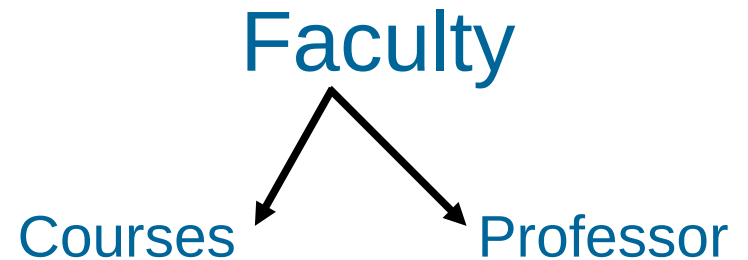
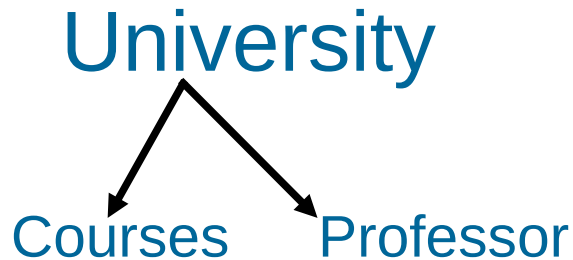
$3\text{grams}(\text{GradCourses}, \text{Courses}) = 0.2$

$\text{Lev}(\text{GradCourses}, \text{Courses}) = 0.42$

$\text{StringMatching}(\text{GradCourses}, \text{Courses}) = 0.31$

$\text{StringMatching}(\text{Professor}, \text{FullProfessor}) = 0.38$

BtreeMatch: Semantic Aspect (3/4)



$\text{StringMatching}(\text{Faculty}, \text{University}) = 0.002$

$\text{Context}(\text{University}) = \{\text{University}, \text{Courses}, \text{Professor}\}$

$\text{Context}(\text{Faculty}) = \{\text{Faculty}, \text{Courses}, \text{Professor}\}$

$\text{CosineMeasure}(\text{Context}(\text{University}), \text{Context}(\text{Faculty})) = 0.37$



BtreeMatch: Semantic Aspect (4/4)

LIRMM

- Acceptable quality of matching
 - better than COMA++ in some scenarios

	Precision	Recall	F-measure
COMA++	1	0.56	0.72
BtreeMatch	0.62	0.89	0.73

- can be tuned
 - to restrict the context
 - to increase the similarity and replacement thresholds



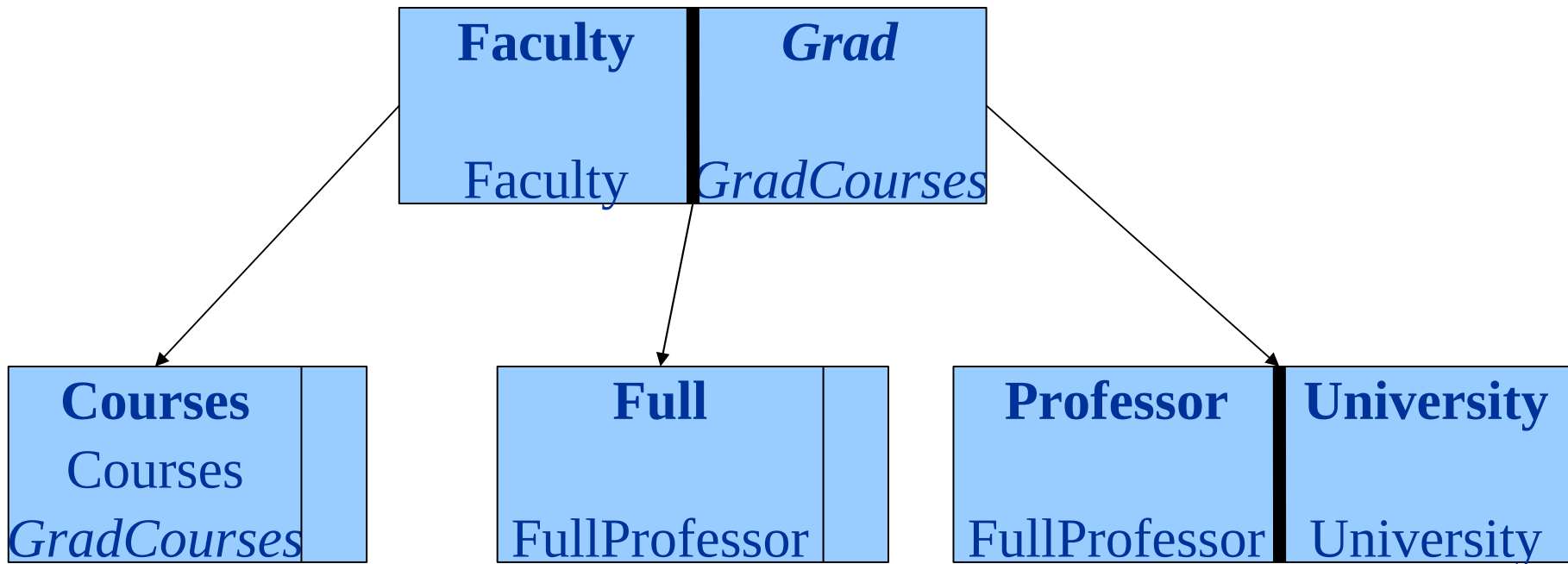
BtreeMatch: Performance Aspect

LIRMM

- B-tree indexing structure to restrict the search space because “*most of similar labels share a common token*”
- Algorithm
 - Labels are divided into tokens
 - Each token is an index in the B-tree with references to all labels containing this token
 - Match search of a label is limited to the labels referenced by the common tokens
 - Else the whole B-tree may be searched using the cosine measure

BtreeMatch: Performance Aspect

LIRMM

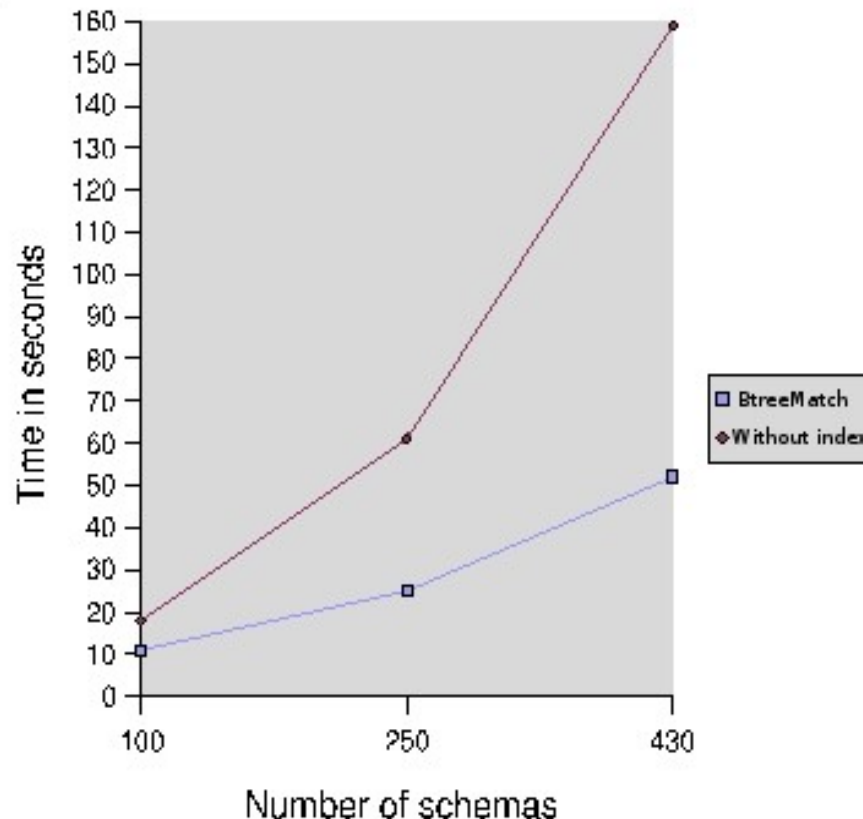


- Searching a match for *GradCourses* involves
 - creation of an index for ***Grad***
 - only evaluating and discovering a similarity between *GradCourses* and *Courses* due to their common token

BtreeMatch: Performance Results

LIRMM

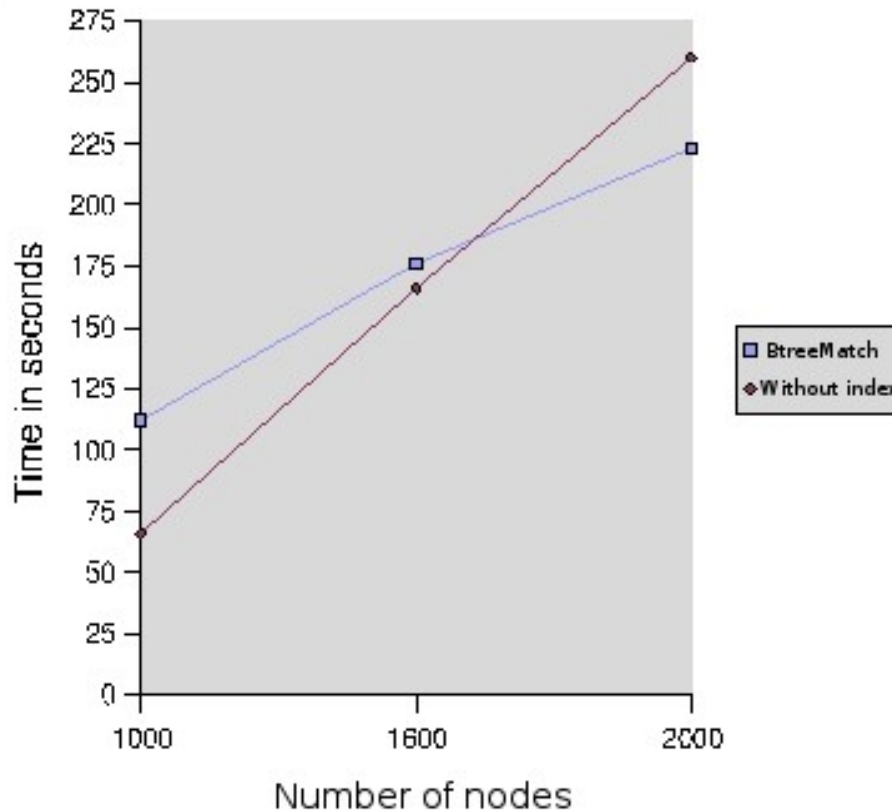
- Comparison of the performance with and without the indexing structure, depending on the number of schemas using XCBL and OASIS schemas



BtreeMatch: Performance Results

LIRMM

- Comparison of the performance with and without the indexing structure, depending on the size of the schemas using XCBL and OASIS schemas





- An automatic schema matching tool that
 - handles many large schemas.
 - provides an acceptable quality of matching.
 - tuning is not automatic

- Discovering complex mappings

- Exploring other index structures (hashtables)