

Benchmark et automatisation du tuning des algorithmes de matching de schémas XML

Encadrement : Nabila BENHARKAT
Rami RIFAIEH

Présenté par : Mohamed BOUKHEBOUZE



Contexte

- Les données des différents systèmes d'information sont hétérogènes.
- Les utilisateurs ont besoin d'une vue intégrée.



Le *matching* de schémas

Contexte

- « *Le matching de schéma est une tâche qui consiste à trouver les correspondances sémantiques entre les éléments de deux schémas* » [Do et Rahm 2002].
- Domaines d'application :
 - Web sémantique,
 - L'intégration de schémas et d'ontologies,
 - L'e-commerce.

Contexte

- Le *matching* a été effectué jusque là manuellement à l'aide d'une interface graphique

Problème

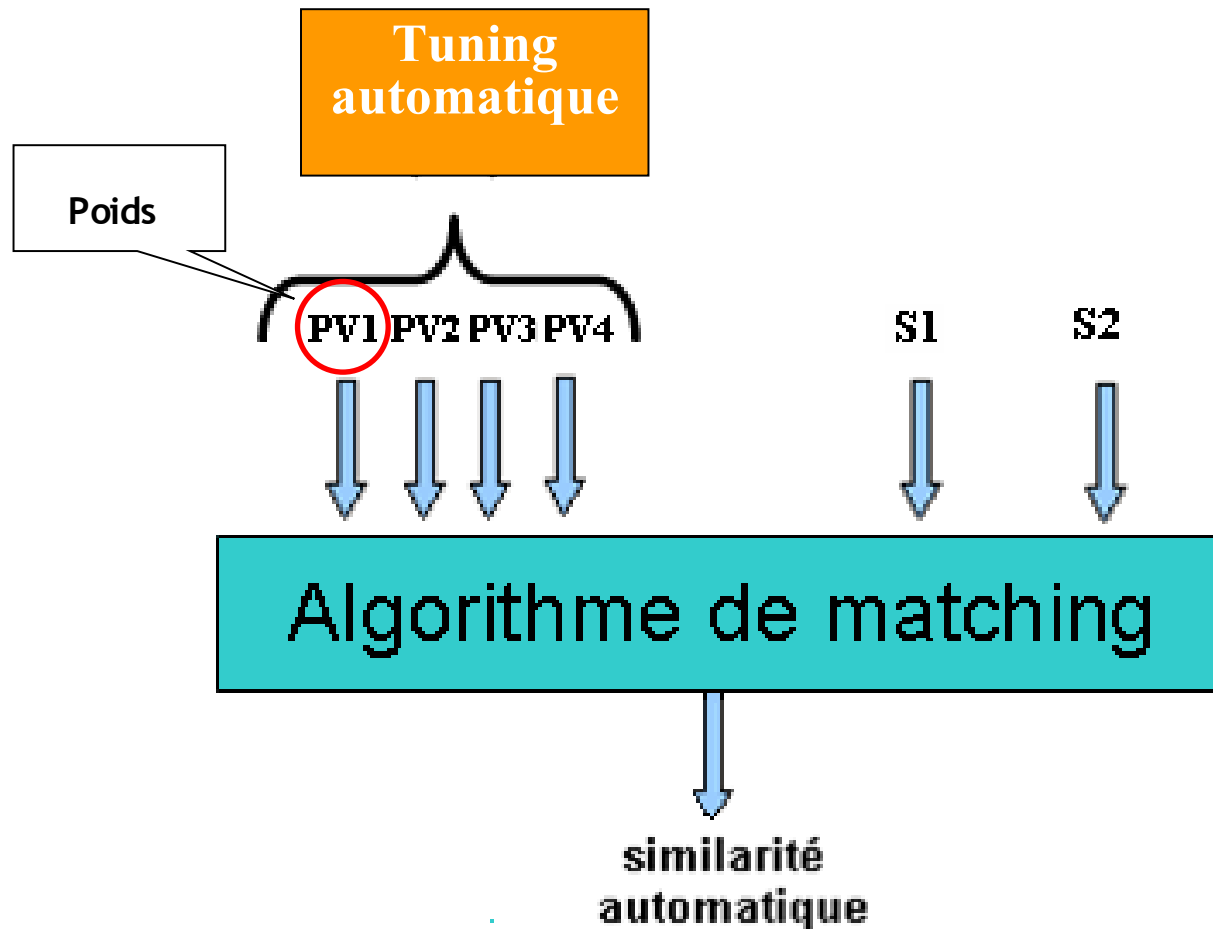
- Consommation du temps
- Intervention des experts humains



Automatiser le processus de *matching*

- Plusieurs algorithmes ont été proposés : COMA [Rahm 02], EXSMAL [Chukmol & al 04], SCIA [Wang 04], Cupid [Madhavan 01], SF [Sergey 01] ET Xmapper [Kurgan 02]

Contexte



Contexte

État de l'art

Notre Approche

*Étapes
benchmark*

*Algorithme
tuning
automatique*

Conclusion

Perspectives

Objectifs

- Déterminer les meilleures valeurs de paramètres qui servent à calculer la similarité structurelle.
- Proposer un benchmark qui examine comment les différents paramètres structurels des algorithmes influent sur les résultats.
- Proposer un algorithme pour automatiser le choix des paramètres structurels optimisés en se basant sur les relations déduites à partir du benchmark.

Contexte

État de l'art

Notre Approche

*Étapes
benchmark*

*Algorithme
tuning
automatique*

Conclusion

Perspectives

Benchmark

- Un benchmark est utile pour comparer l'efficacité de différents systèmes matching
 - ➔ Très peu de benchmarks de matching proposés
- Des évaluations sur les algorithmes de matching qui sont toujours individuelles

Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Évaluation de COMA

- Cinq schémas XML avec des caractéristiques différentes
- Dix tâches de matching sont définies, pour chaque tâche deux différents schémas
- 12312 séries de données sont générées. On compare les résultats obtenus avec les résultats obtenus avec les autres schémas. On indique le nombre de vraies correspondances parmi celles retournées et le nombre de vraies correspondances parmi celles retournées.
- Une fois les résultats obtenus, on calcule la précision, le Recall et le Overall. La précision, le Recall et le Overall sont calculés

indique le nombre de vraies correspondances parmi celles retournées

indique le nombre de vraies correspondances parmi celles retournées

$$\text{Recall} * \left(2 - \frac{1}{\text{precision}}\right)$$

Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Tuning automatique

- Le tuning est la tâche qui consiste à ajuster les paramètres des différents matchers
- ➔ Il existe peu de travaux proposés pour tuning automatique :
- [Berkovsky 2005](#) : les algorithmes génétiques pour le tuning automatique.
 - [Sayyadian 2005](#) : un système complètement automatique (eTuner)

Contexte

État de l'art

Notre Approche

**Étapes
benchmark**

**Algorithme
tuning
automatique**

Conclusion

Perspectives

Les algorithmes génétiques pour le tuning automatique

- + L'utilisation des AG permet d'avoir une solution générique au problème de tuning automatique
- L'efficacité des algorithmes génétiques dépend du choix de la fonction objective

Contexte

État de l'art

Notre Approche

**Étapes
benchmark**

**Algorithme
tuning
automatique**

Conclusion

Perspectives

eTuner

+ Le eTuner a de bonnes performances

- Passe par deux étapes

- Perturbation des schémas
- Exécution du tuning

 coût considérable en temps.

Contexte

État de l'art

Notre Approche

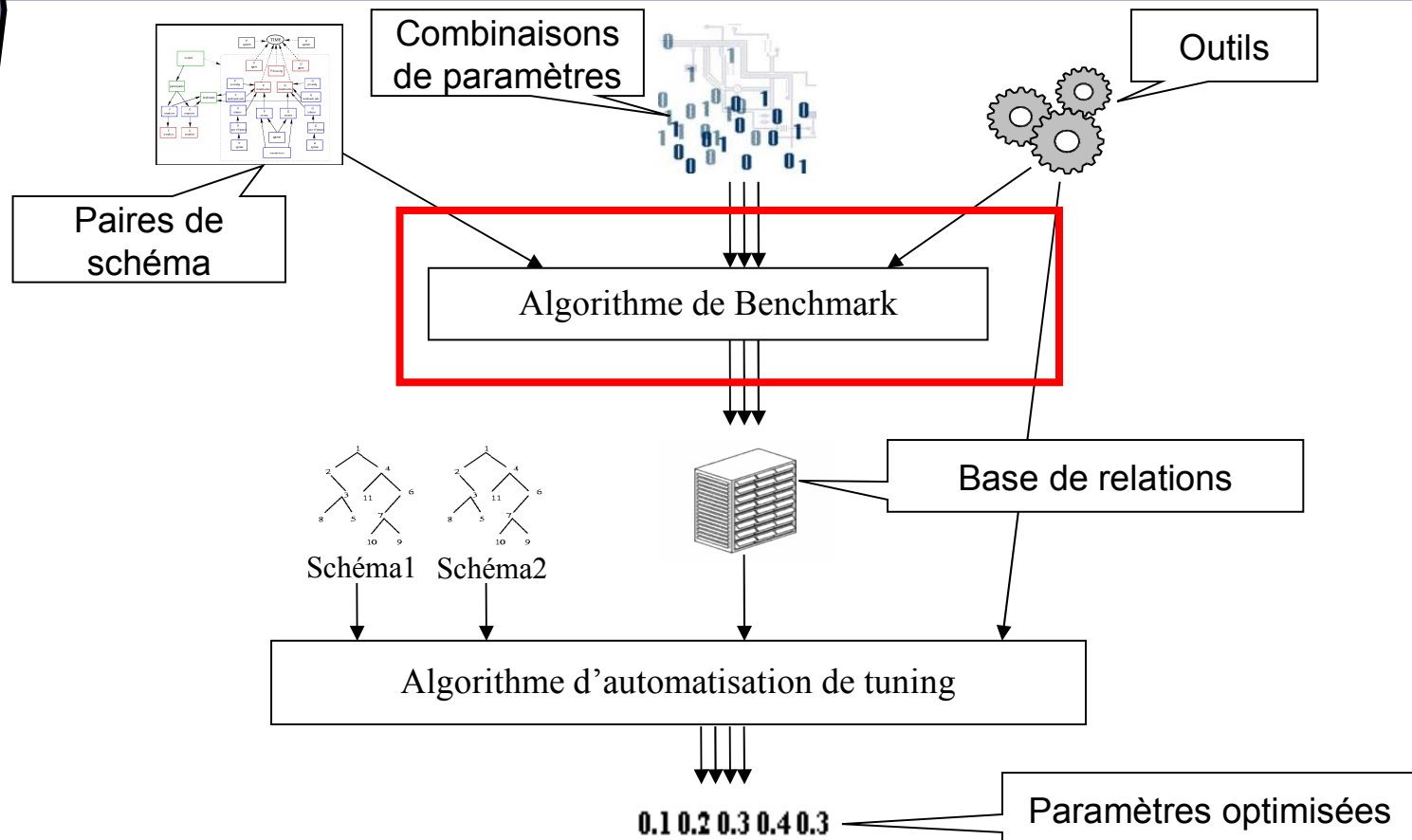
Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Notre approche



Contexte

État de l'art

Notre Approche

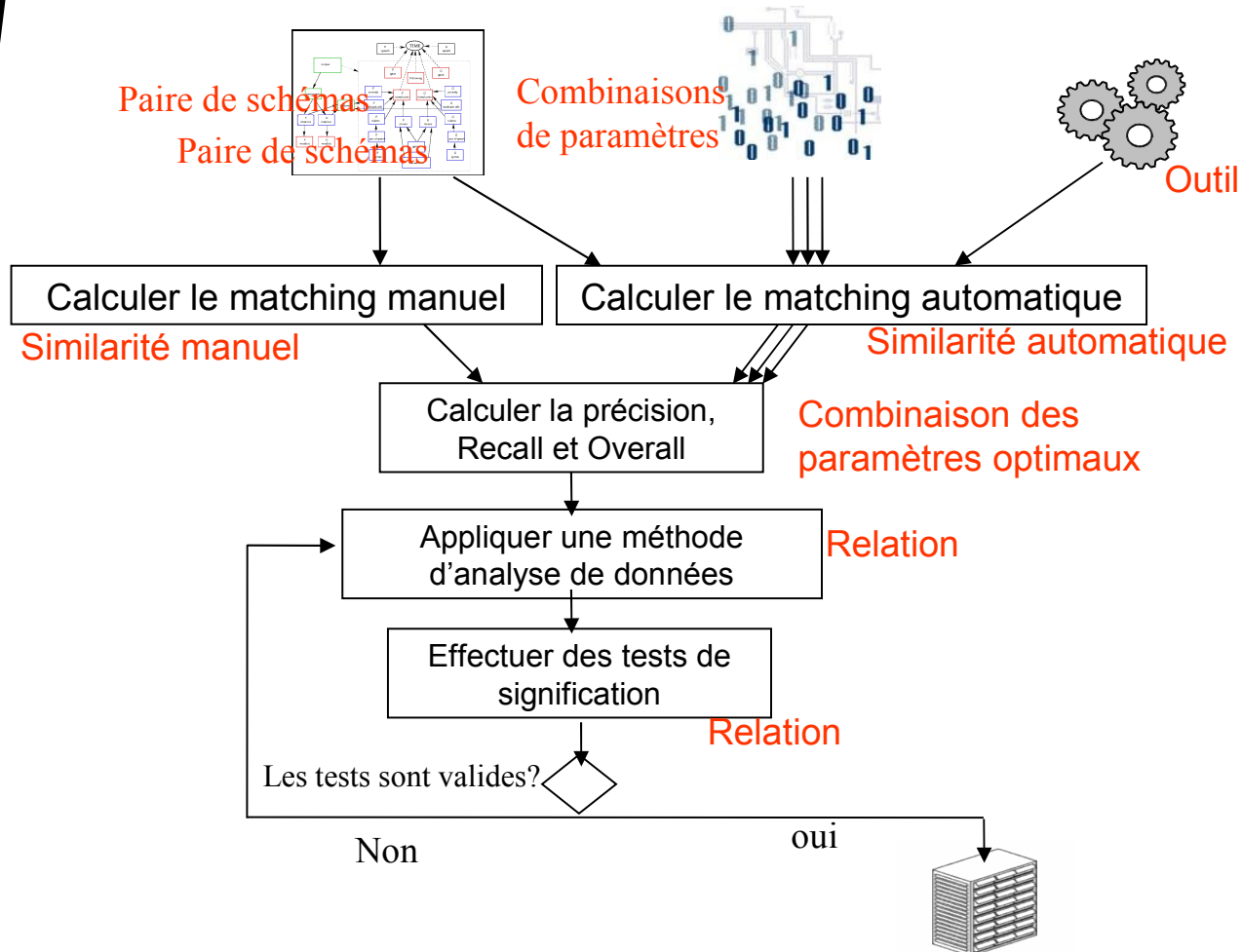
Étapes benchmark

Algorithme tuning automatique

Conclusion

Perspectives

Étape de benchmark



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Schémas de tests

- 20 schémas différents utilisés :

$$N, P \in]0,80] \text{ et } D \in [2,7]$$

Où N : Nombre d'éléments

P : Nombre de chemins

D : La profondeur

➔ $C_{20}^2 = 190$ combinaisons de paires

de schémas possibles

Étape de benchmark

Algorithme de matching choisis

Contexte

État de l'art

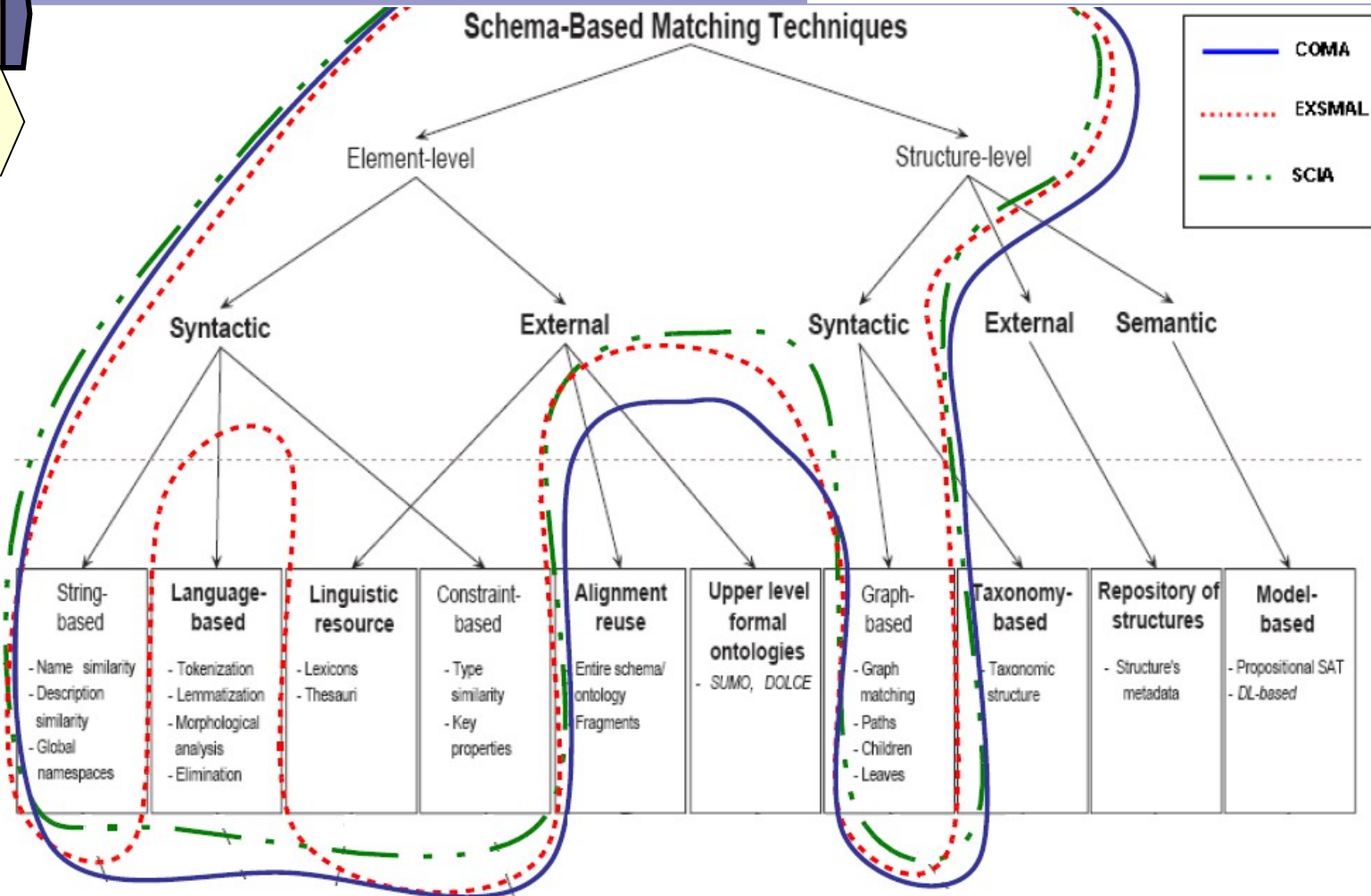
Notre Approche

Étapes benchmark

Algorithme tuning automatique

Conclusion

Perspectives



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Paramètres structurels étudiés

- EXSMAL : 4 Paramètres
combinaisons possible → 20
- COMA : 4 Paramètres
combinaisons possible → 625
- SCIA : 3 Paramètres
combinaisons possible → 125

Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Nombres de séries de tests

■ EXSMAL : $190 \times 20 = 3800$

■ COMA : 118

■ SCIA : 237

Nombre de
combinaisons de
paires de schémas

Nombre de
combinaisons de
paramètres



146300

Nombre de séries
de tests

Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Algorithme de benchmark

Pour chaque($O_i \in EO$) faire //(nous avons choisi 3 outils :EXSMAL ,
COMA et SCIA)

Pour chaque ($(S_i, S_j) \in ES$) faire // (190 paire de schémas possibles)

Pour chaque (combinaison de paramètres $\in EP$) faire

// (les combinaisons possibles de paramètres avec une incrémentation
delta)

// Calculer le matching avec la paire de schémas;

//Calculer la précision, le Recall et le Overall;

(PR, RE, OV) Matching (Si, Sj, EC1);

Fin_Pour

//Trouver le modèle de régression

//Appliquer les tests de signification

Si les tests sont valides alors Sauvegarder le modèle

Sinon Appliquer d'autres méthodes d'analyse de données

Fin_Si

Fin_Pour

Fin_Pour

Contexte

État de l'art

Notre Approche

Étapes
benchmark

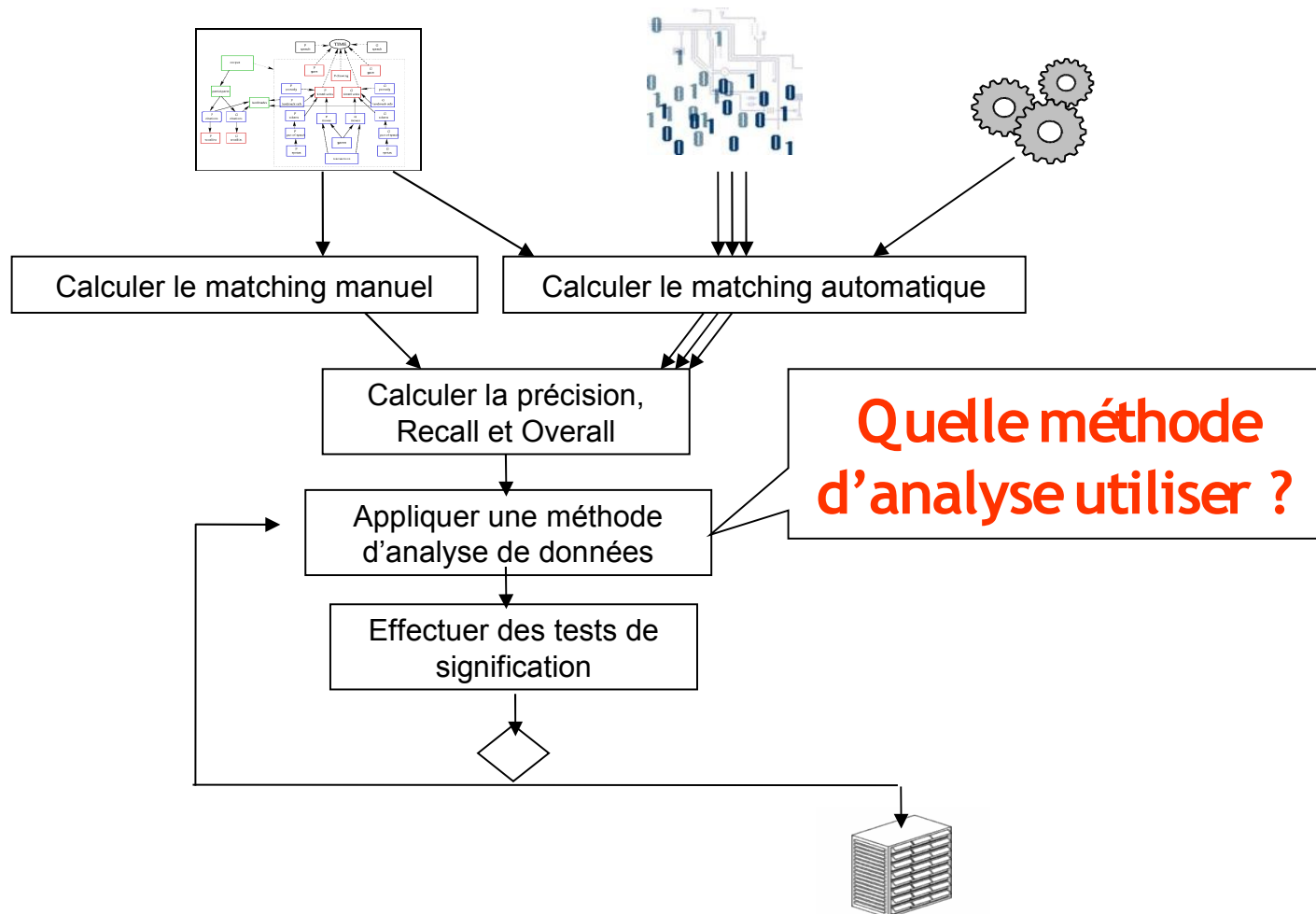
Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Benchmark pour EXSMAL



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

comment juger que la
relation est
significative ?

Étape de benchmark

Analyse de données

- Quelle méthode d'analyse de données utilisée?

- Les méthodes :

- Les ACP

- La régression multiple

Réduire le nombre de variables à étudier

Calculer l'équation de l'hyperplan

$$Y = a_0 + a_1 x_1 + a_2 x_2 + e$$

variable expl

vari

variable explicative

Erreur de mesure

Contexte

État de l'art

Notre Approche

Étapes
benchmark


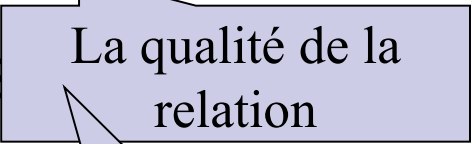
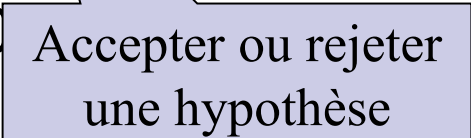
Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Tests de significations

- Analyse des résidus 
 - Le coefficient de détermination
 - Les tests d'hypothèses (test) 

- [Gérald Baillargeon, Livre : Méthodes statistiques, les éditions SMG, 2004 p.273-40]

Contexte

État de l'art

Notre Approche

Étapes
benchmark

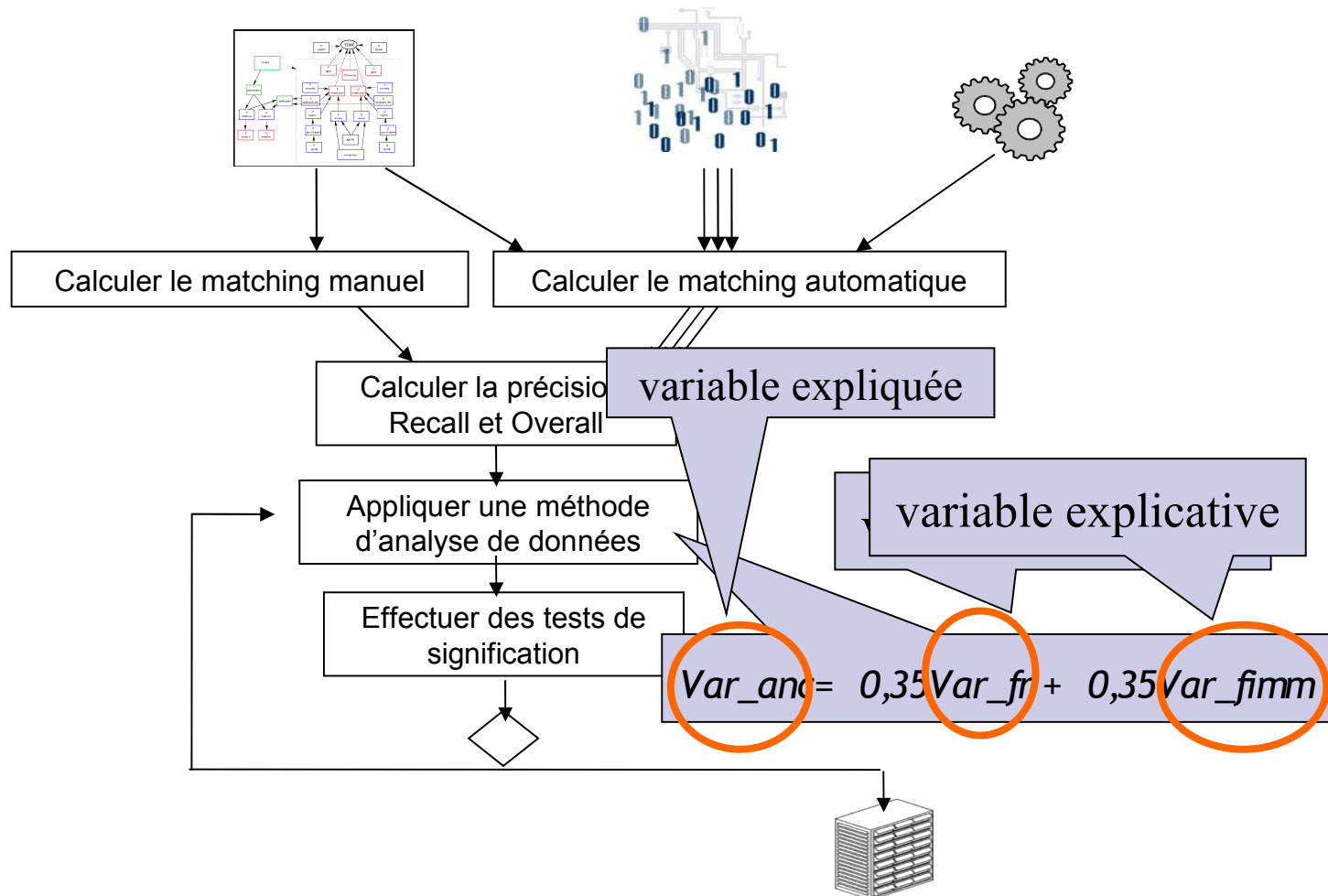
Algorithme
tuning
automatique

Conclusion

Perspectives

Étape de benchmark

Benchmark pour EXSMAL



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Benchmark pour EXSMAL

- coefficient de détermination $R^2 = 0,79$.

- Residus normalisés

Bonne qualité

- L'hypothèse statistique
H0: le *coeff_anc* est linéairement indépendant
des *coeff_fr* et *coeff_fimm*.

Erreurs acceptables

D'après le test de Fisher

Hypothèses rejetées

Contexte

État de l'art

Notre Approche

Étapes
benchmark

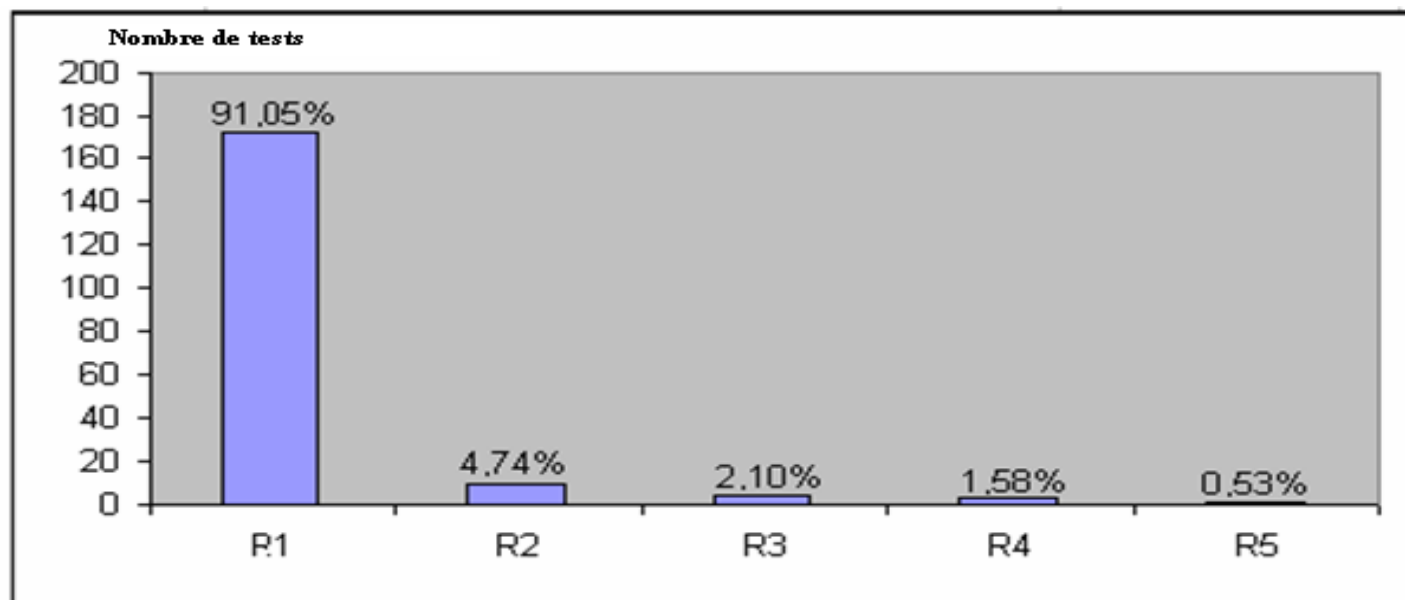
Algorithme
tuning
automatique

Conclusion

Perspectives

Benchmark pour EXSMAL

| Relation | La relation mathématique | N et P | D | Nombre de tests |
|----------|---|--------|-------|-----------------|
| R1 | $\text{Coeff_anc} = 0.35\text{Coeff_fr} + 0.35\text{Coeff_fimm}$ | [0,80] | [2,7] | 173 |
| R2 | $\text{Coeff_anc} = 0.33\text{Coeff_fr} + 0.33\text{Coeff_fimm}$ | [0,80] | [2,7] | 9 |
| R3 | $\text{Coeff_anc} = 0.25\text{Coeff_fr} + 0.04\text{Coeff_fimm}$ | [0,80] | [2,7] | 4 |
| R4 | $\text{Coeff_anc} = 0.13\text{Coeff_fr} + 0.27\text{Coeff_fimm}$ | [0,80] | [2,7] | 3 |
| R5 | $\text{Coeff_anc} = 2.2\text{Coeff_fr} + 0.45\text{Coeff_fimm}$ | [0,80] | [2,7] | 1 |



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Benchmark pour COMA

- $P_Parents = a0 + a1P_Children + a2P_Leaves$
- $P_Siblings = 10$

| Relation | La relation mathématique | N et P | D | Nombre de tests |
|----------|--|---------|-------|-----------------|
| R1 | $P_Parents = 0,6 P_Children + 0,6 P_Leaves$ |]0,40[| [2,4] | 94 |
| R2 | $P_Parents = 0,48 P_Children + 0,48 P_Leaves$ | [40,80] | [2,4] | 27 |
| R3 | $P_Parents = 0,11 P_Children + 0,11 P_Leaves$ |]0,40[| [5,7] | 51 |
| R4 | $P_Parents = 0,07 P_Children + 0,07 P_Leaves$ | [40,80] | [5,7] | 18 |

Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Benchmark pour SCIA

- $PATH_WEIGHT = a0 + a1GRAPH_WEIGHT + a2STRUCT_WEIGHT$

| Relation | La relation mathématique | N et P | D | nombre de test |
|----------|---|---------|-------|----------------|
| R1 | $PATH_WEIGHT = -0,23$ $GRAPH_WEIGHT - 0,62$ $STRUCT_WEIGHT + 0,51$ |]0,40[| [2,4] | 94 |
| R2 | $PATH_WEIGHT = -0,5$ $GRAPH_WEIGHT - 0,5$ $STRUCT_WEIGHT + 0,5$ | [40,80] | [2,4] | 27 |
| R3 | $PATH_WEIGHT = -0,25$ $GRAPH_WEIGHT - 0,75$ $STRUCT_WEIGHT + 0,55$ |]0,40[| [5,7] | 51 |
| R4 | $PATH_WEIGHT = -0,49$ $GRAPH_WEIGHT - 0,33$ $STRUCT_WEIGHT + 0,52$ | [40,80] | [5,7] | 18 |

Contexte

État de l'art

Notre Approche

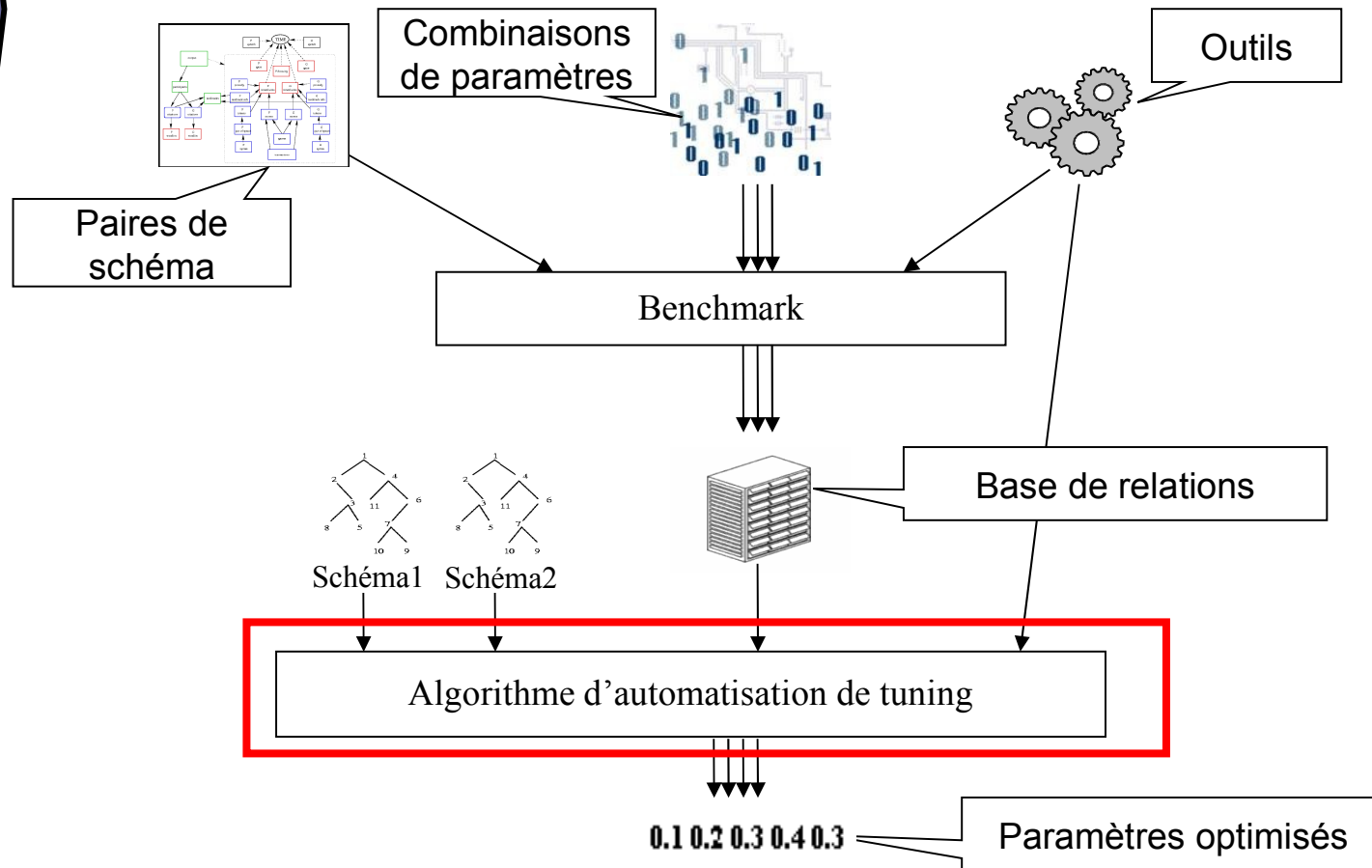
Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Notre approche



Contexte

État de l'art

Notre Approche

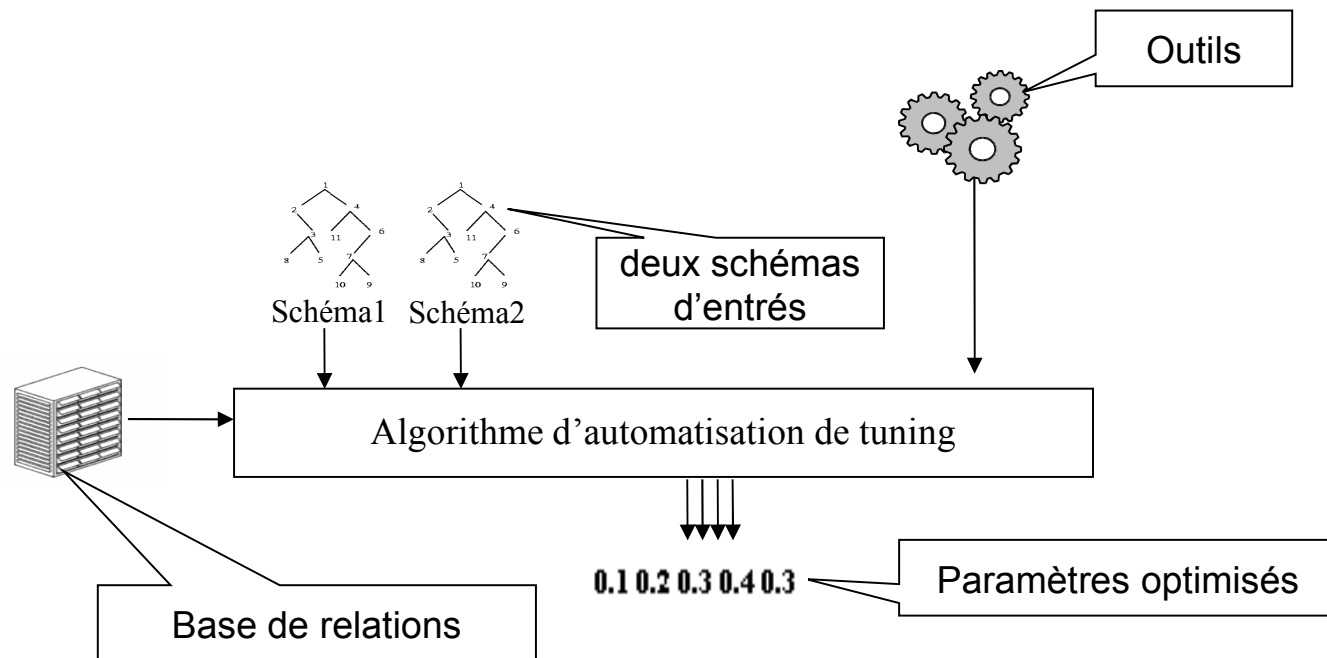
Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Algorithme de tuning



Contexte

État de l'art

Notre Approche

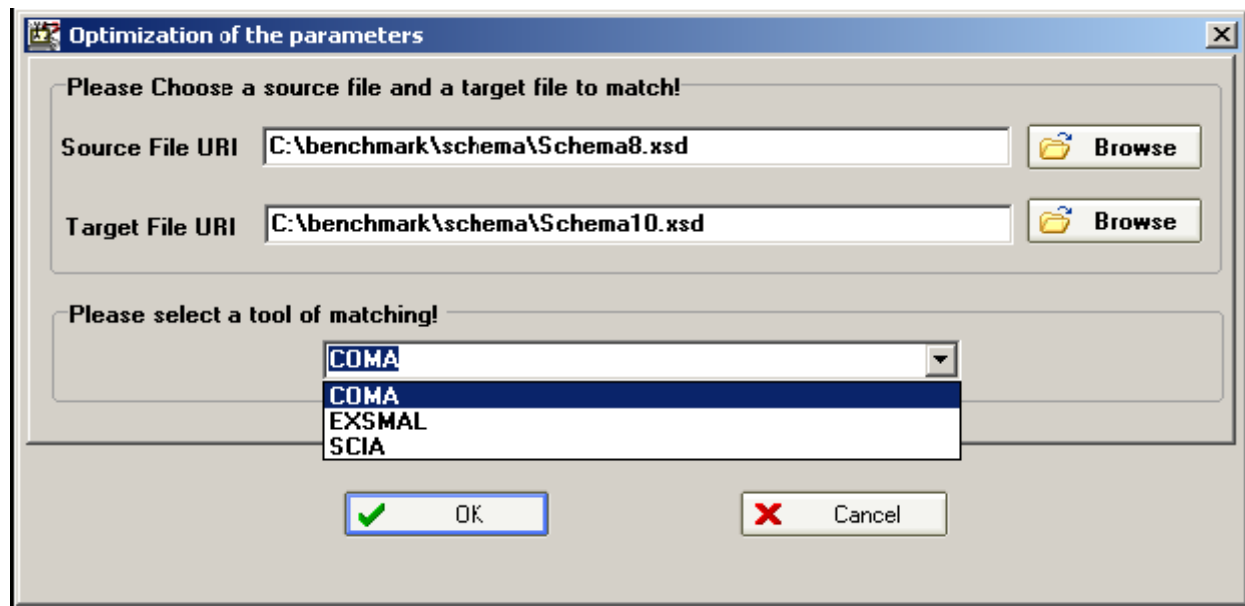
Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Algorithme de tuning



Contexte

État de l'art

Notre Approche

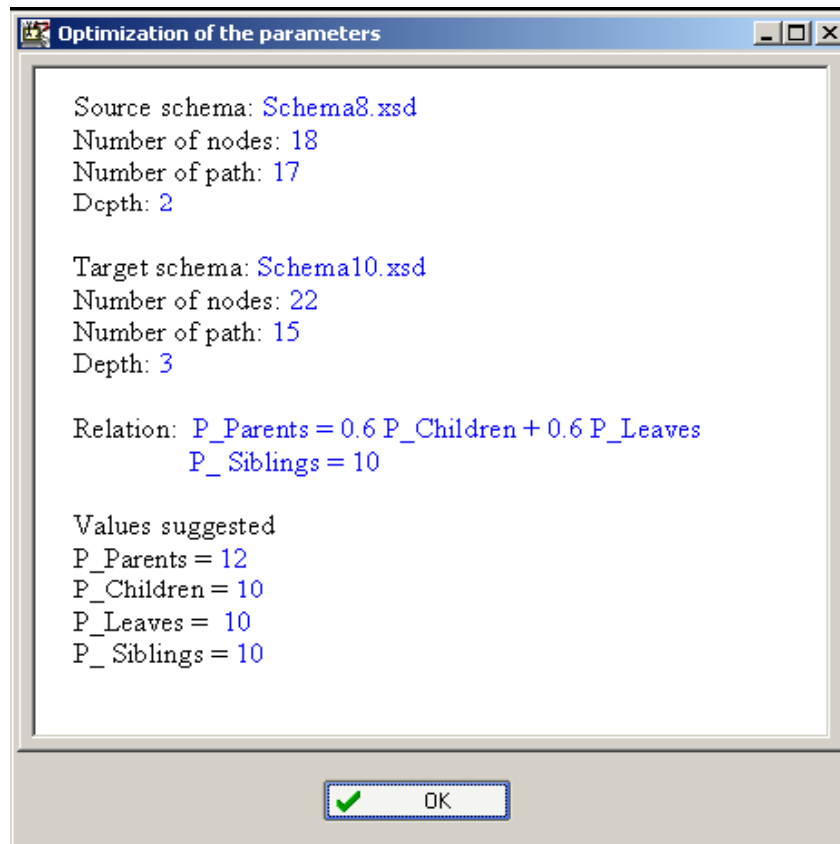
Étapes
benchmark

Algorithme
tuning
automatique

Conclusion

Perspectives

Algorithme de tuning



Contexte

État de l'art

Notre Approche

Étapes
benchmark

Algorithme
tuning
automatique

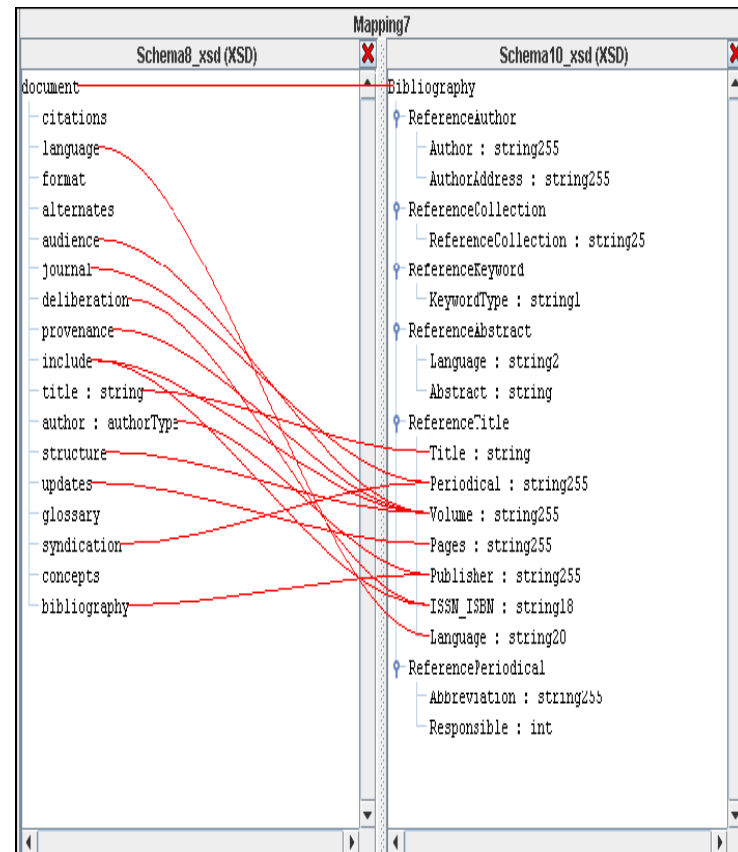
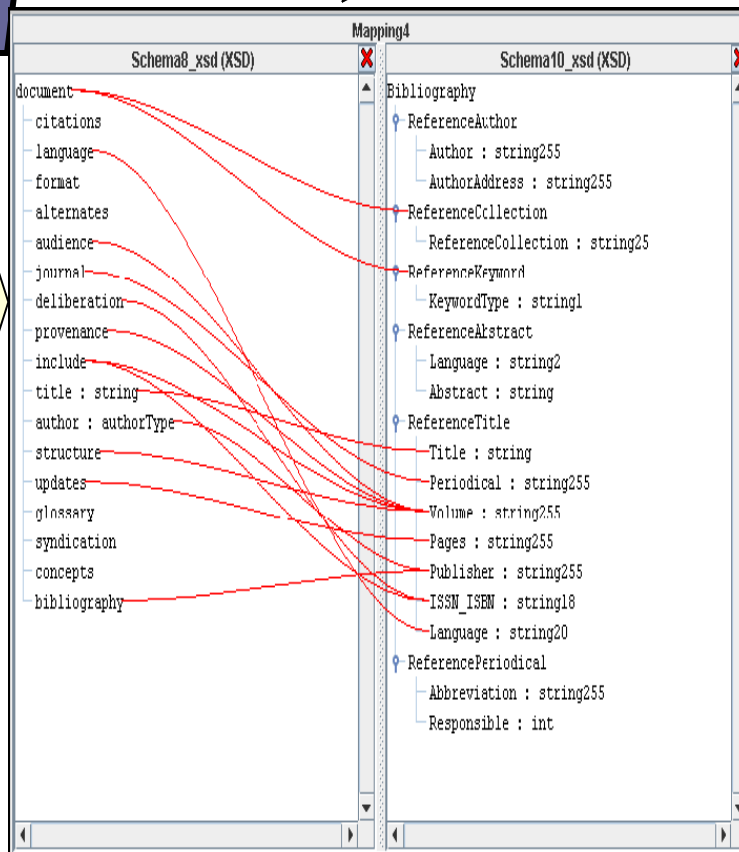
Conclusion

Perspectives

Algorithme de Tuning

Précision = 0,23

Précision = 0,48



Le matching avant l'exécution de l'outil

Le matching après l'exécution de l'outil

Contexte

État de l'art

Notre Approche

*Étapes
benchmark*

*Algorithme
tuning
automatique*

Conclusion

Perspectives

Conclusion

- Solution statistique
- Benchmark évolutif
- Algorithme de tuning évolutif
- Notre proposition par rapport à eTuner?

Contexte

État de l'art

Notre Approche

**Étapes
benchmark**

**Algorithme
tuning
automatique**

Conclusion

Perspectives

Perspectives

- Tuning des paramètres descriptifs
- Réaliser le benchmark avec un nombre de schémas plus important, et d'autres outils
- Intégrer cette approche dans la plateforme ASMADE

*Je vous remercie de votre
attention*

