

# Réponse à des requêtes conjonctives en présence d'ontologies

Décidabilité, complexité et algorithmes

Michaël Thomazo

Encadré par Jean-François Baget et Marie-Laure Mugnier  
Université Montpellier 2

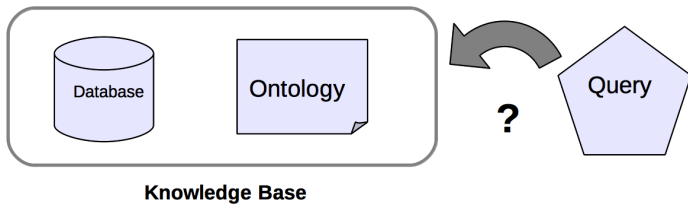
30 juin 2014

RJCIA, Rouen

## Vers une interrogation sémantique

- ▶ interrogation “intelligente” du Web (ou d’un sous-ensemble : Wikipedia par exemple) ;
- ▶ traitement sémantique de données bio-médicales.

# Vue synthétique du problème



# À la croisée de domaines

- ▶ représentation des connaissances et raisonnement
- ▶ base de données

# Représentation des connaissances et raisonnement

- ▶ représenter des connaissances de manière formelle → ontologies
- ▶ ontologies : théorie logique représentant les connaissances d'un domaine
- ▶ principale approche pour exprimer des ontologies : les logiques de description
- ▶ raisonnement centré sur l'ontologie elle-même (historiquement)

## Bases de données – un exemple

<b>memberOf</b>	
Alice	FencingClub
Bob	FencingClub
Craig	GardensFriends

<b>SportClub</b>
FencingClub

<b>employeeOf</b>	
Bob	LIRMM
Craig	_x1

$q_1(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y)$

## Bases de données – un exemple

memberOf	
Alice	FencingClub
Bob	FencingClub
Craig	GardensFriends

SportClub
FencingClub

employeeOf	
Bob	LIRMM
Craig	_x1

$q_1(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y)$       {Alice, Bob}

## Bases de données – un exemple

memberOf	
Alice	FencingClub
Bob	FencingClub
Craig	GardensFriends

SportClub
FencingClub

employeeOf	
Bob	LIRMM
Craig	_x1

$q_1(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y)$       {Alice, Bob}

$q_2(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y) \wedge \text{hasSocialSecurityNumber}(x, z)$



## Bases de données – un exemple

memberOf	
Alice	FencingClub
Bob	FencingClub
Craig	GardensFriends

SportClub
FencingClub

employeeOf	
Bob	LIRMM
Craig	_x1

$q_1(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y)$       {Alice, Bob}

$q_2(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y) \wedge \text{hasSocialSecurityNumber}(x, z)$        $\emptyset$

## Bases de données – un exemple

memberOf	
Alice	FencingClub
Bob	FencingClub
Craig	GardensFriends

SportClub
FencingClub

employeeOf	
Bob	LIRMM
Craig	_x1

$q_1(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y)$       {Alice, Bob}

$q_2(x) : \neg \text{memberOf}(x, y) \wedge \text{SportClub}(y) \wedge \text{hasSocialSecurityNumber}(x, z)$        $\emptyset$

“Tout employé a un numéro de sécurité sociale.”

# Un problème à l'intersection des deux domaines

## Ontology Based Query Answering (OBQA)

Interroger des données en prenant en compte une ontologie.

# Outline

## Panorama

Règles existentielles

Approches usuelles

Limites

## Une nouvelle classe décidable

Bounded Treewidth Sets

Greedy Bounded Treewidth Sets

## Les outils pour un algorithme optimal

## Résumé et travaux futurs

# Panorama

## Règles existentielles (1)

$$\forall \mathbf{x} \forall \mathbf{y} (B[\mathbf{x}, \mathbf{y}] \rightarrow \exists \mathbf{z} H[\mathbf{y}, \mathbf{z}]),$$

Une ontologie est un ensemble de règles existentielles.

Même forme logique que :

- ▶ TGDs [Abiteboul et al., 94]
- ▶ Datalog+/- [Calì et al., 09]
- ▶ règles de graphes conceptuels [Sowa 84, Mugnier and Salvat 96]

## Règles existentielles (2)

- ▶ “value invention”
- ▶ pas de restrictions d'arité
- ▶ corps et tête de règles de forme quelconque
- ▶ généralise les logiques de description légères, couramment utilisées pour OBQA

## Un exemple de règles existentielles

“Chaque employé a un numéro de sécurité sociale.”

$$\forall x \forall y ( \text{employeeOf}(x, y) \rightarrow \exists z \text{ hasSocialSecurityNumber}(x, z) )$$

Le numéro de sécurité sociale n'est pas précisé – juste son existence.



## Sur l'exemple

- ▶ données :  
 $\text{memberOf}(\text{Alice}, \text{SportsCenter}) \wedge \dots$
- ▶ ontologie :  
 $\{\forall x \forall y (\text{employeeOf}(x, y) \rightarrow \exists z \text{hasSSN}(x, z)), \dots\}$
- ▶ requête :  
 $\exists x \exists y \exists z \text{memberOf}(x, y) \wedge \text{SportClub}(y) \wedge \text{hasSSN}(x, z)$

Des requêtes non-booléennes peuvent aussi être considérées.

## Idée centrale

Se ramener à une “simple” requête SQL à évaluer sur une base de données.

## Application de règles

$\forall x \forall y \text{ employeeOf}(x,y) \rightarrow \exists z \text{ hasSocialSecurityNumber}(x, z)$

memberOf( Alice, SportsCenter)  $\wedge$   
memberOf( Bob, SportsCenter)  $\wedge$   
memberOf( Craig, SwimmingCenter)  $\wedge \dots \wedge$   
employeeOf( Bob,LIRMM)  $\wedge$  employeeOf( Craig,\_x1)

- ▶  $x \rightarrow \text{Bob}, y \rightarrow \text{LIRMM}$  :  
création de  $\text{hasSocialSecurityNumber}(\text{Bob}, z_1)$

## Application de règles

$\forall x \forall y \text{ employeeOf}(x,y) \rightarrow \exists z \text{ hasSocialSecurityNumber}(x, z)$

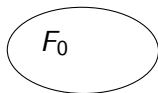
memberOf( Alice, SportsCenter)  $\wedge$   
memberOf( Bob, SportsCenter)  $\wedge$   
memberOf( Craig, SwimmingCenter)  $\wedge \dots \wedge$   
employeeOf( Bob,LIRMM)  $\wedge$  employeeOf( Craig, \_x1)

- ▶  $x \rightarrow \text{Bob}, y \rightarrow \text{LIRMM}$  :  
création de  $\text{hasSocialSecurityNumber}(\text{Bob}, z_1)$
- ▶  $x \rightarrow \text{Craig}, y \rightarrow \_x1$  :  
création de  $\text{hasSocialSecurityNumber}(\text{Craig}, z_2)$

$z_1$  et  $z_2$  sont existentiellement quantifiées.

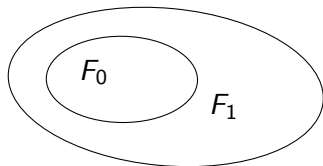
# Modèle canonique (matérialisation)

Chase, marche avant :  $F$  et  $\mathcal{R}$  en entrée



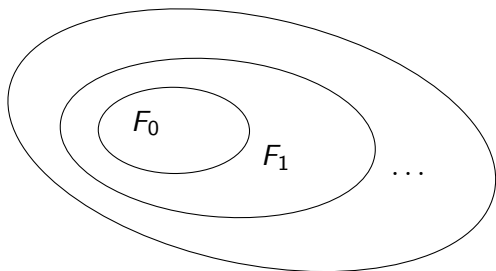
# Modèle canonique (matérialisation)

Chase, marche avant :  $F$  et  $\mathcal{R}$  en entrée



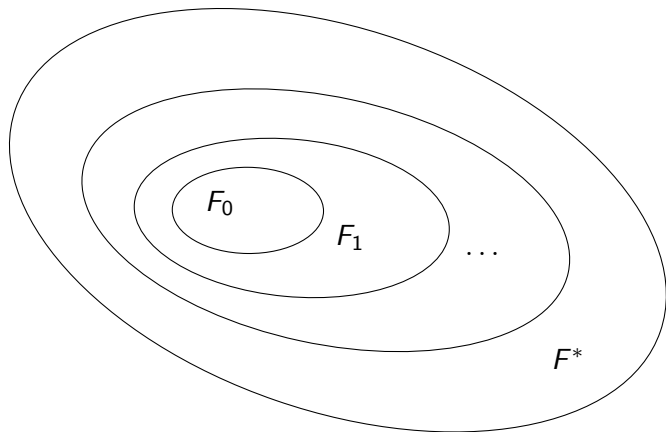
# Modèle canonique (matérialisation)

Chase, marche avant :  $F$  et  $\mathcal{R}$  en entrée



# Modèle canonique (matérialisation)

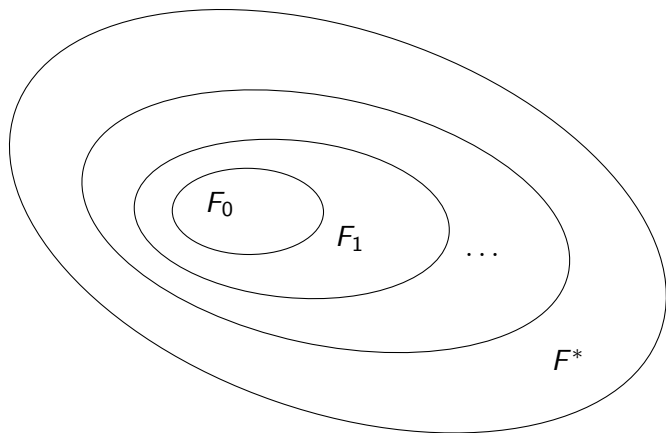
Chase, marche avant :  $F$  et  $\mathcal{R}$  en entrée





# Modèle canonique (matérialisation)

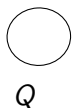
Chase, marche avant :  $F$  et  $\mathcal{R}$  en entrée



**Property :**  $\forall Q (F, \mathcal{R} \models Q \Leftrightarrow F^* \models Q)$

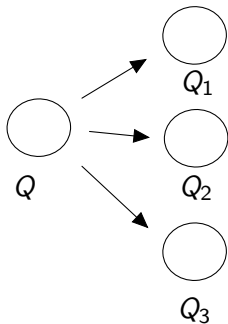
# Sans matérialisation ?

Réécriture, marche arrière :  $Q$  and  $\mathcal{R}$  en entrée ([Calvanese et al., 2005] entre autres)



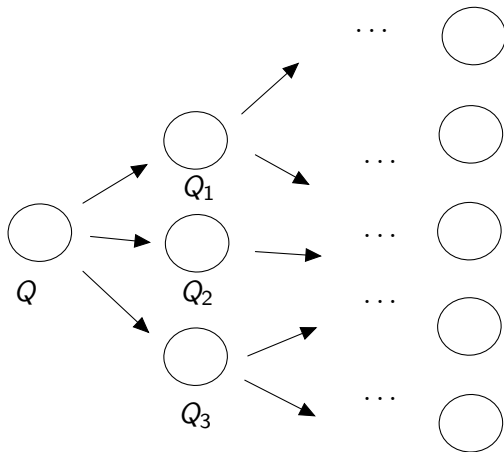
# Sans matérialisation ?

Réécriture, marche arrière :  $Q$  and  $\mathcal{R}$  en entrée ([Calvanese et al., 2005] entre autres)



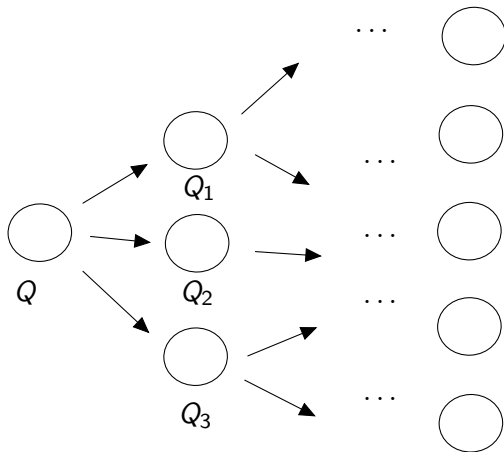
# Sans matérialisation ?

Réécriture, marche arrière :  $Q$  and  $\mathcal{R}$  en entrée ([Calvanese et al., 2005] entre autres)



# Sans matérialisation ?

Réécriture, marche arrière :  $Q$  and  $\mathcal{R}$  en entrée ([Calvanese et al., 2005] entre autres)



**Propriété :**  $\forall F (F, \mathcal{R} \models Q \Leftrightarrow \exists Q_i : F \models Q_i)$

# Indécidabilité

- ▶ le modèle canonique peut être infini :  
 $\forall x ( \text{Human}(x) \rightarrow \text{parent}(y, x) \wedge \text{Human}(y) )$
- ▶ une réécriture de taille finie peut ne pas exister :  
 $\forall x \forall y ( \text{parent}(x, y) \wedge \text{Human}(x) \rightarrow \text{Human}(y) )$

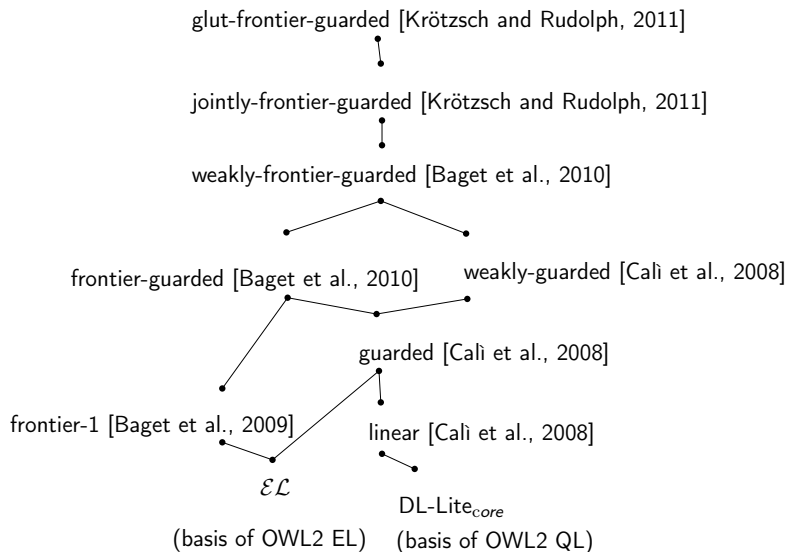
**Indécidabilité** [Beeri and Vardi 1981], entre autres

Répondre à des requêtes conjonctives en présence de règles existentielles n'est pas décidable.

# Principaux critères de décidabilité

- ▶ quand le modèle canonique est équivalent à un fait fini ;
- ▶ quand il existe une réécriture finie ;
- ▶ quand le modèle canonique est “semblable” à un arbre (*bts* : bounded treewidth set).

## Classes concrètes *bts* connues





# Contributions

- ▶ définition d'une nouvelle classe décidable ;
- ▶ construction d'un algorithme optimal dans le pire des cas pour cette classe ;
- ▶ adaptation de l'algorithme à des sous-classes ;
- ▶ design, prototype et première évaluation d'un algorithme générique de réécriture.

# Contributions

- ▶ définition d'une nouvelle classe décidable ;
- ▶ construction d'un algorithme optimal dans le pire des cas pour cette classe ;
- ▶ adaptation de l'algorithme à des sous-classes ;
- ▶ design, prototype et première évaluation d'un algorithme générique de réécriture.

# Outline

## Panorama

Règles existentielles

Approches usuelles

Limites

## Une nouvelle classe décidable

Bounded Treewidth Sets

Greedy Bounded Treewidth Sets

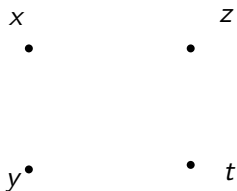
## Les outils pour un algorithme optimal

## Résumé et travaux futurs

# Une nouvelle classe décidable

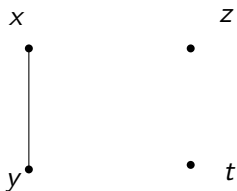
# Graph primal d'un fait

Grphe primal de  $\exists x \exists y \exists z \exists t \ r(x, y) \wedge p(x, z, t) \wedge r(y, z)$



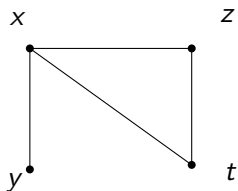
# Graph primal d'un fait

Grphe primal de  $\exists x \exists y \exists z \exists t \ r(x, y) \wedge p(x, z, t) \wedge r(y, z)$



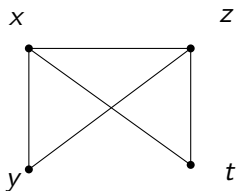
# Graph primal d'un fait

Grphe primal de  $\exists x \exists y \exists z \exists t \ r(x, y) \wedge p(x, z, t) \wedge r(y, z)$



# Graph primal d'un fait

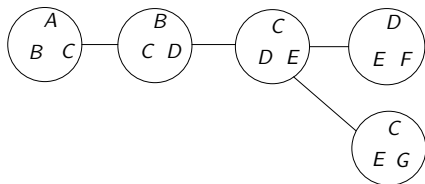
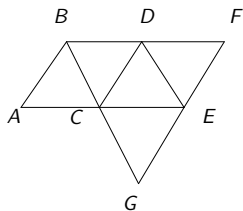
Grphe primal de  $\exists x \exists y \exists z \exists t \ r(x, y) \wedge p(x, z, t) \wedge r(y, z)$





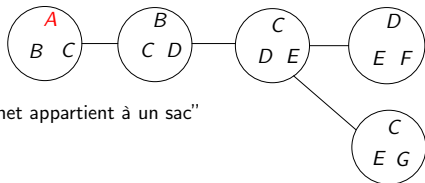
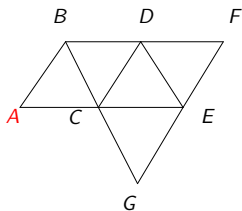
# Décomposition et largeur arborescentes (1)

Robertson and Seymour, 1984



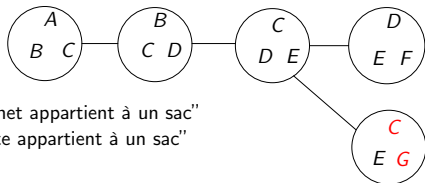
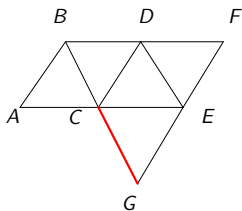
# Décomposition et largeur arborescentes (1)

Robertson and Seymour, 1984



# Décomposition et largeur arborescentes (1)

Robertson and Seymour, 1984

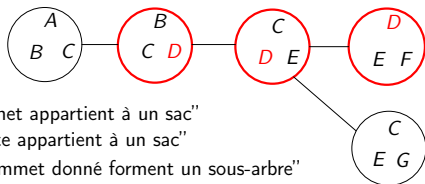
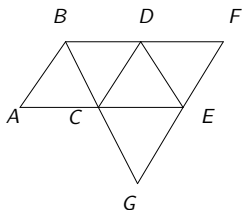


"Chaque sommet appartient à un sac"

"Chaque arête appartient à un sac"

# Décomposition et largeur arborescentes (1)

Robertson and Seymour, 1984



"Chaque sommet appartient à un sac"

"Chaque arête appartient à un sac"

"Les sacs contenant un sommet donné forment un sous-arbre"

## Décomposition et largeur arborescentes (2)

### Treewidth

La largeur d'une décomposition arborescente est le nombre maximal de sommets dans un sac moins un. La largeur arborescente (treewidth) d'un graphe est la plus petite largeur d'une de ses décompositions arborescentes.

### Treewidth d'un arbre

La treewidth d'un arbre est 1.

A quel point ressemble un graphe à un arbre ?

# Décidabilité de *bts*

Calí et al. 2008, Baget et. al 2009

## Treewidth d'un fait

La treewidth d'un fait est la treewidth de son graph primal.

## Bounded treewidth sets (*bts*)

Un ensemble de règles est à treewidth bornée si pour tout fait  $F$ , la treewidth de son modèle canonique est bornée.

## OBQA avec *bts* est décidable

Cela provient d'un résultat de Courcelle sur la logique monadique du second ordre (1990).

Pas de manière imposée d'obtenir une décomposition – pas d'algorithmes connus

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011

$$r(a, b, c), p(b)$$

$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

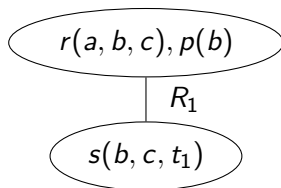
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

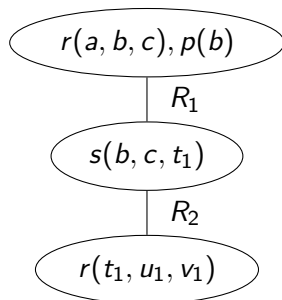
$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$



# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

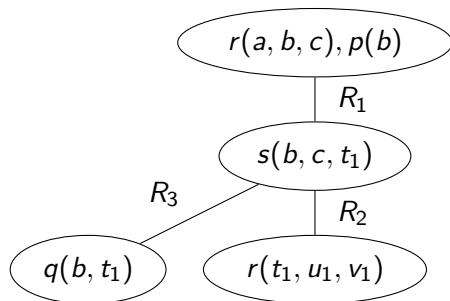
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

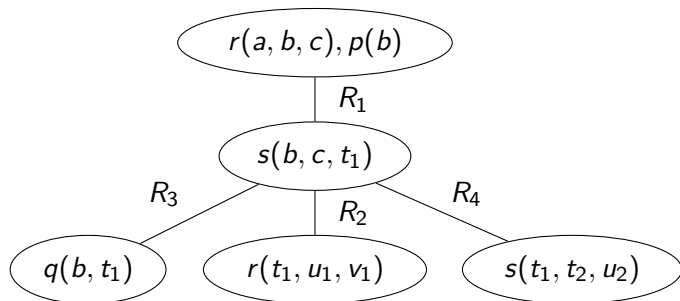
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

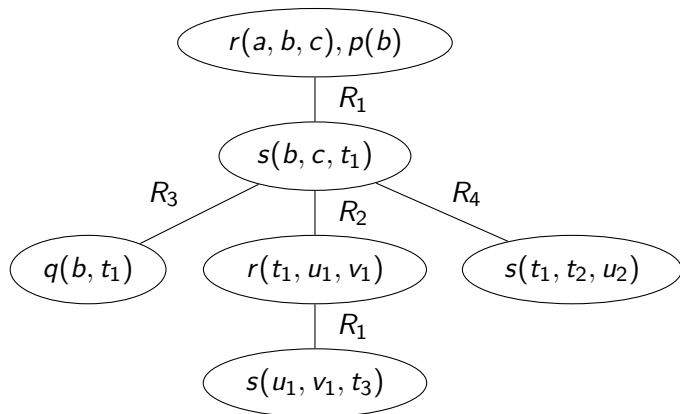
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

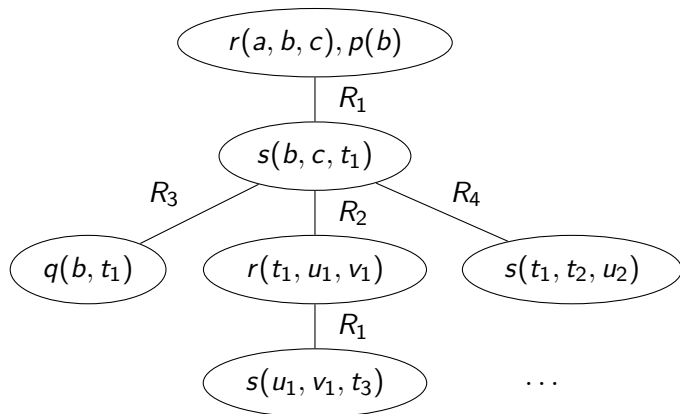
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

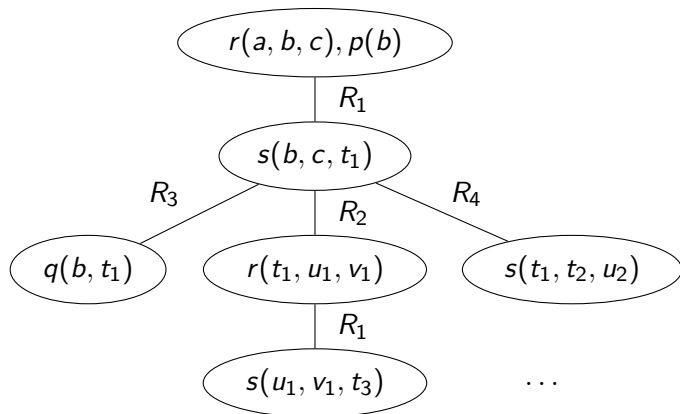
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

# Décomposition arborescente gloutonne

Baget, Mugnier, Rudolph, T., 2011



$$R_1 = r(x, y, z) \rightarrow s(y, z, t)$$

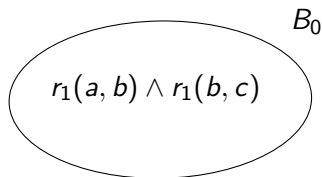
$$R_3 = s(x, y, z) \wedge p(x) \rightarrow q(x, z)$$

$$R_2 = s(x, y, z) \rightarrow r(z, u, v)$$

$$R_4 = q(x, z) \rightarrow s(z, t, u)$$

Ultimate derivation tree

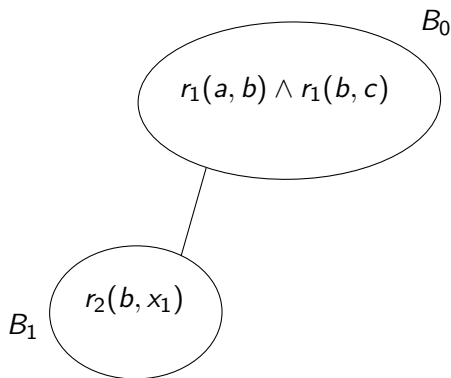
## Un exemple de règles qui est *bts* mais pas *gbts*



$$R_1 = r_1(x, y) \rightarrow r_2(y, z)$$

$$R_2 = r_2(x, y) \wedge r_1(x, z) \wedge r_2(z, t) \rightarrow r_3(y, t)$$

## Un exemple de règles qui est *bts* mais pas *gbts*

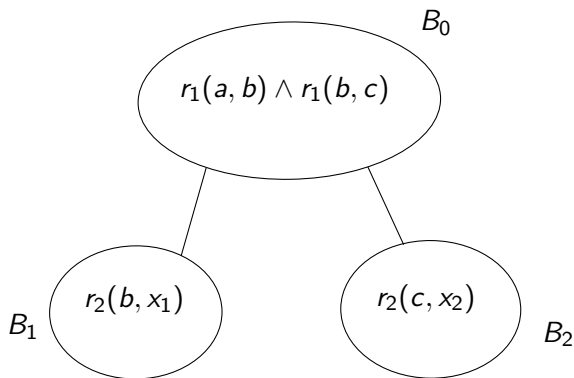


$$R_1 = r_1(x, y) \rightarrow r_2(y, z)$$

$$R_2 = r_2(x, y) \wedge r_1(x, z) \wedge r_2(z, t) \rightarrow r_3(y, t)$$



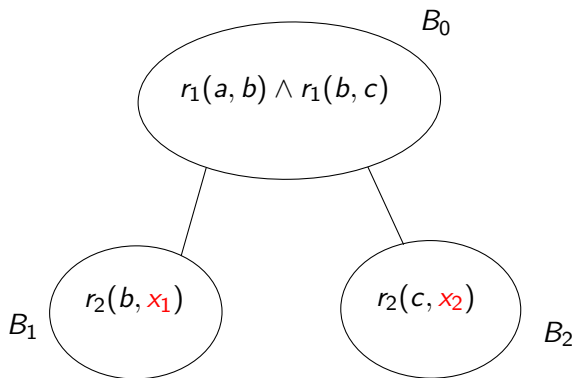
## Un exemple de règles qui est *bts* mais pas *gbts*



$$R_1 = r_1(x, y) \rightarrow r_2(y, z)$$

$$R_2 = r_2(x, y) \wedge r_1(x, z) \wedge r_2(z, t) \rightarrow r_3(y, t)$$

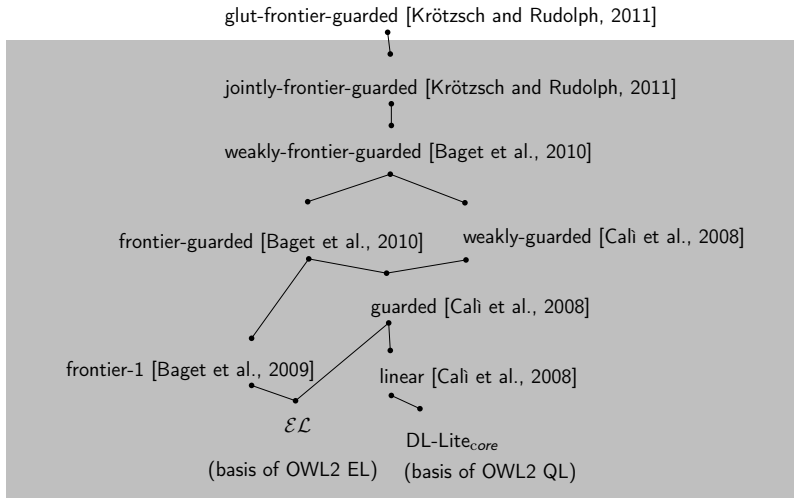
## Un exemple de règles qui est *bts* mais pas *gbts*



$$R_1 = r_1(x, y) \rightarrow r_2(y, z)$$

$$R_2 = r_2(x, y) \wedge r_1(x, z) \wedge r_2(z, t) \rightarrow r_3(y, t)$$

# gbts Classes



frontier-1 : la frontière est de taille 1

guarded : un atome du corps de la règle contient toutes ses variables

# Idées principales de l'algorithme pour *gbts*

## Représenter le modèle canonique de manière finie.

Principaux outils et étapes :

- ▶ définir une relation d'équivalence convenable sur les sacs de la décomposition arborescente du modèle canonique (*patterns* and *patterns abstraits*);
- ▶ calculer cette relation d'équivalence sans construire cette décomposition arborescente (*evolution* and *creation rules*);
- ▶ construire la représentation elle-même.

# Outline

## Panorama

Règles existentielles

Approches usuelles

Limites

## Une nouvelle classe décidable

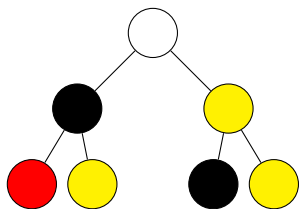
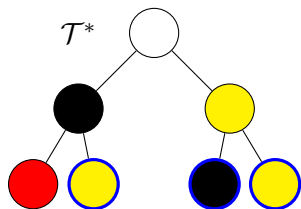
Bounded Treewidth Sets

Greedy Bounded Treewidth Sets

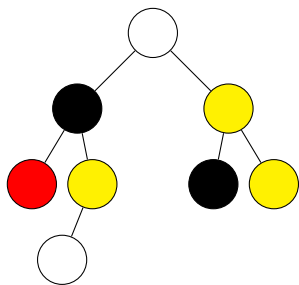
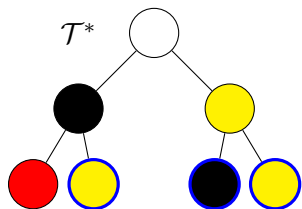
## Les outils pour un algorithme optimal

## Résumé et travaux futurs

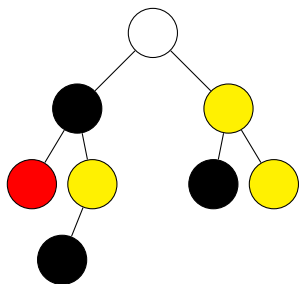
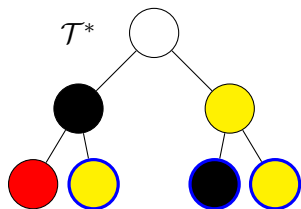
## Représentation finie de la décomposition arborescente



## Représentation finie de la décomposition arborescente

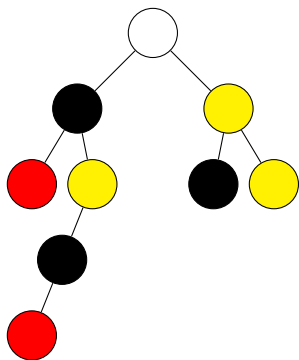
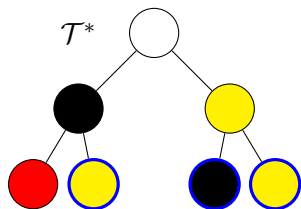


## Représentation finie de la décomposition arborescente

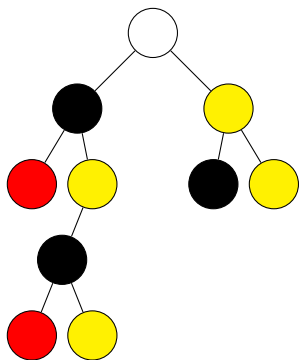
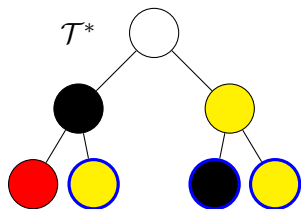




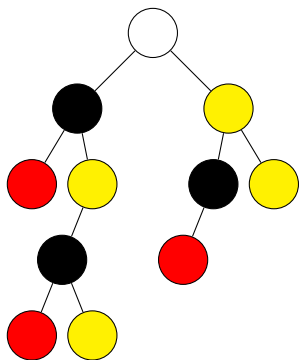
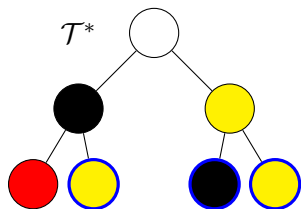
## Représentation finie de la décomposition arborescente



## Représentation finie de la décomposition arborescente



## Représentation finie de la décomposition arborescente



...

# Patterns et patterns abstraits

- ▶ relation d'équivalence se sert de *patterns*
- ▶ les sacs d'un arbre de dérivation sont étiquetés par des patterns
- ▶ intuitivement : cela représente comment les corps de règles sont envoyés dans ce sac

## Equivalence

Deux sacs avec même pattern abstrait sont à la racine d'arbres isomorphes dans la décomposition arborescente du modèle canonique

# Calcul des patterns abstraits

- ▶ **but** : calcul des patterns abstraits
- ▶ **outil** : règles d'évolution et de création

# Complexité

Classe	Combinée	Données
<i>gbts</i>	2ExpTime-c [T. et al. 12]	ExpTime-c [Baget et al. 11]
<b>j-fg</b>	2ExpTime-c [Krötzsch et al. 11]	ExpTime-c [Krötzsch et al. 11]
<b>wfg</b>	2ExpTime-c [Baget et al. 11]	ExpTime-c [Baget et al. 11]
<b>fg</b>	2ExpTime-c [Baget et al. 11]	PTime-c [Baget et al. 11]
<b>fr1</b>	2ExpTime-c [Baget et al. 11]	PTime-c [Baget et al. 11]
<b>wg</b>	2ExpTime-c [Cali et al. 08]	ExpTime-c [Cali et al. 08]
<b>guarded</b>	2ExpTime-c [Cali et al. 08]	PTime-c [Cali et al. 08]
<b>linear</b>	PSpace-c [Cali et al. 10]	in AC <sub>0</sub> [Cali et al. 09]

## Résumé

- ▶ définition d'une nouvelle classe décidable et expressive pour OCQA [Baget, Mugnier, Rudolph, T., IJCAI'11]
- ▶ construction d'un algorithme optimal dans le pire des cas pour cette classe  
[T., Baget, Mugnier, Rudolph, KR'12]
- ▶ adaptation de cet algorithme à des sous-classes [T., Baget, Mugnier, Rudolph, KR'12] et [T., DL'12]
- ▶ design, implementation d'un prototype and et première évaluation d'un algorithme de réécriture [T., IJCAI'13]

## Travaux futurs

- ▶ implémentation de l'algorithme pour *gbts*
- ▶ extension de l'approche
- ▶ réécriture au delà des formules du premier-ordre

## Remerciements

- ▶ à Jean-François et Marie-Laure pour leur encadrement au top ;
- ▶ à GraphIK pour l'environnement ;
- ▶ au jury (Georg Gottlob, Carsten Lutz, Christophe Paul, Marie-Christine Rousset) ;
- ▶ au jury du prix de thèse AFIA ;
- ▶ à la DAAD et à la JSPS ;
- ▶ et à vous, pour votre attention !