

Interrogation de données médiatisée par des connaissances : le cas des key-value stores

Olivier RODRIGUEZ

<olivier.rodriquez@inria.fr>

Lundi 17 décembre 2018

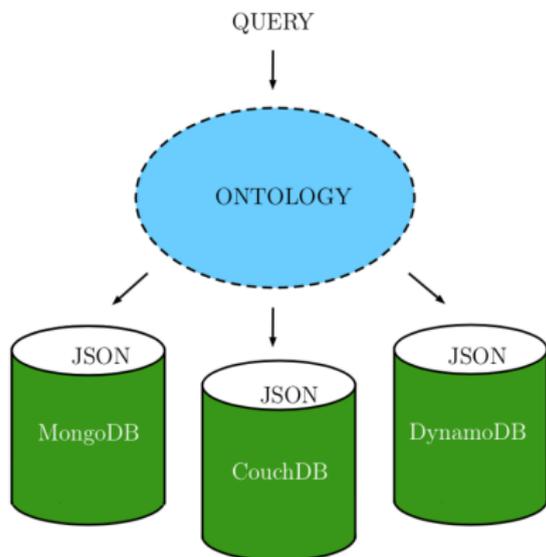
Données arborescentes

Pourquoi s'y intéresser ?

- Données très utilisées
- Systèmes de stockages efficaces & récents
- Spécialisation des graphes

Cadre de travail

Ontology-Mediated Query Answering (OMQA)



Avantages du cadre :

- Enrichir le vocabulaire d'interrogation
- Gérer de l'information incomplète
- Vue unifiée pour des données hétérogènes

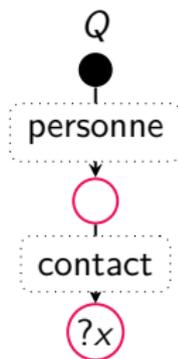
Problématiques de recherche :

- *Interrogation en présence d'ontologie*
- Interrogation fédérée

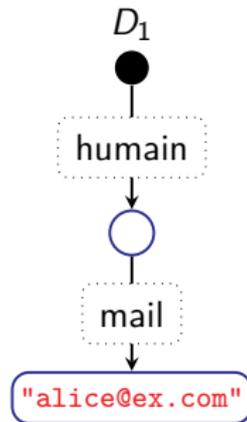
Interrogation en présence d'ontologie

Matérialisation

Requête



Base de données



Règles

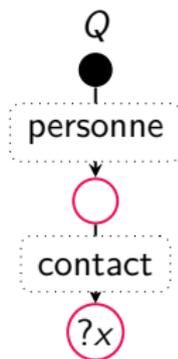
$\text{humain}(X,Y) \longrightarrow \text{personne}(X,Y)$

$\text{mail}(X,Y) \longrightarrow \text{contact}(X,Y)$

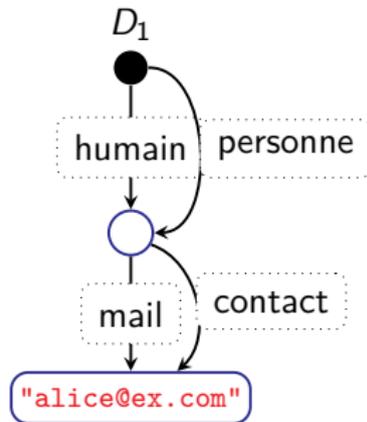
Interrogation en présence d'ontologie

Matérialisation

Requête



Base de données



Règles

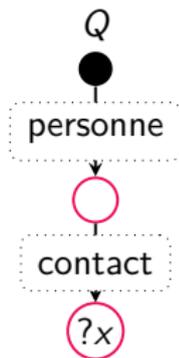
$\text{humain}(X,Y) \longrightarrow \text{personne}(X,Y)$

$\text{mail}(X,Y) \longrightarrow \text{contact}(X,Y)$

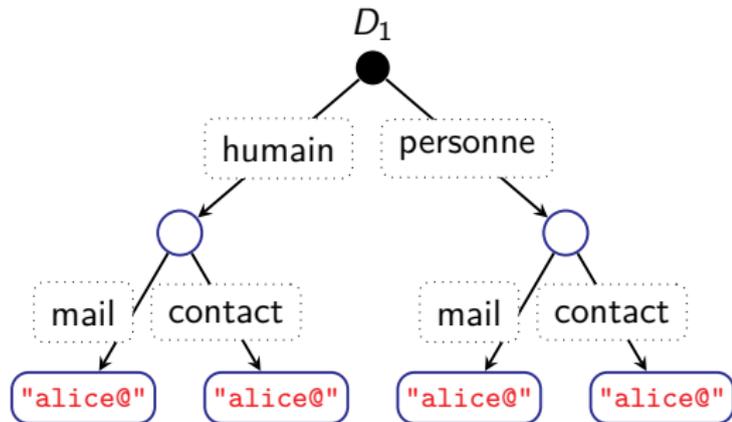
Interrogation en présence d'ontologie

Matérialisation

Requête



Base de données



Règles

$\text{humain}(X,Y) \longrightarrow \text{personne}(X,Y)$

$\text{mail}(X,Y) \longrightarrow \text{contact}(X,Y)$

Interrogation en présence d'ontologie

Matérialisation

Problèmes

- Explosion de l'espace de stockage
- Synchronisation avec les bases sources

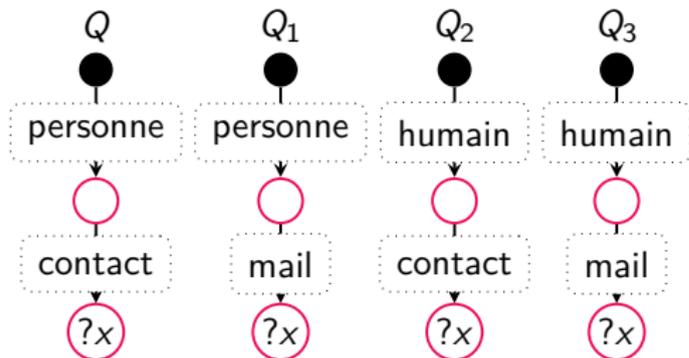
Axe de recherche

Réécriture des requêtes.

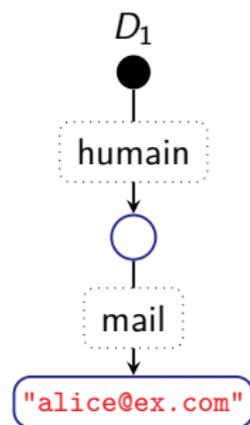
Interrogation en présence d'ontologie

Réécriture

Réécritures



Base de données



Règles

$\text{humain}(X,Y) \longrightarrow \text{personne}(X,Y)$

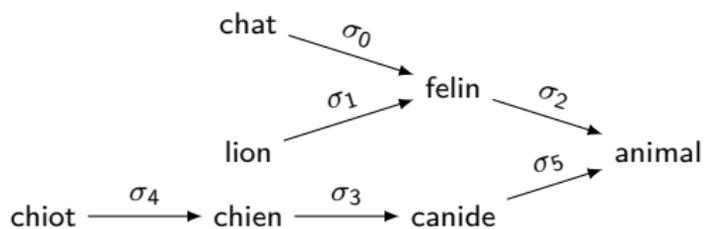
$\text{mail}(X,Y) \longrightarrow \text{contact}(X,Y)$

Règles [MRU16]

$$\begin{array}{l} \text{humain}(X,Y) \longrightarrow \text{personne}(X,Y) \\ \text{personne}(X,Y) \longrightarrow \text{nom}(X,Z) \end{array} \implies \begin{array}{l} \text{humain} \longrightarrow \text{personne} \\ \text{personne} \longrightarrow \exists \text{nom} \end{array}$$

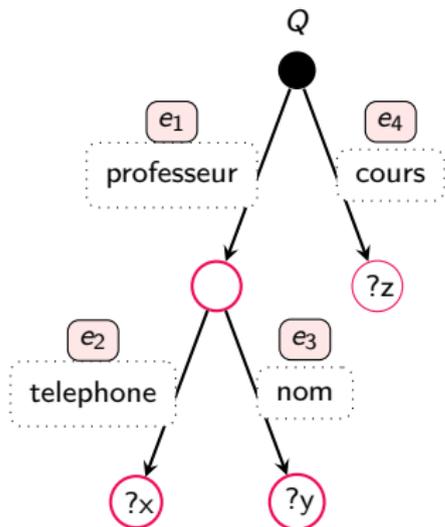
Graphe de dépendances des règles

$$\Sigma = \left\{ \begin{array}{ll} \sigma_0 : \text{chat} & \rightarrow \text{felin} \\ \sigma_1 : \text{lion} & \rightarrow \text{felin} \\ \sigma_2 : \text{felin} & \rightarrow \text{animal} \\ \sigma_3 : \text{chien} & \rightarrow \text{canide} \\ \sigma_4 : \text{chiot} & \rightarrow \text{chien} \\ \sigma_5 : \text{canide} & \rightarrow \text{animal} \end{array} \right\}$$



$$\begin{aligned} \text{app}(\text{animal}, \Sigma) &= \Sigma \\ \text{app}(\text{felin}, \Sigma) &= \{\sigma_0, \sigma_1\} \\ \text{app}(\text{canide}, \Sigma) &= \{\sigma_3, \sigma_4\} \\ \text{app}(\text{chat}, \Sigma) &= \emptyset \end{aligned}$$

Contextes de codages



$$\Sigma = \left\{ \begin{array}{ll} \sigma_1 : \textit{personne} & \rightarrow \textit{professeur} \\ \sigma_2 : \textit{humain} & \rightarrow \textit{personne} \\ \sigma_3 : \textit{mail} & \rightarrow \exists \textit{telephone} \\ \sigma_4 : \textit{UE} & \rightarrow \textit{cours} \end{array} \right\}$$

$$C_c = \left\{ \begin{array}{l} Cc_1 = \{(0, \textit{professeur}), (1, \textit{personne}), (2, \textit{humain})\} \\ Cc_2 = \{(0, \textit{telephone}), (1, \textit{mail})\} \\ Cc_3 = \{(0, \textit{nom})\} \\ Cc_4 = \{(0, \textit{cours}), (1, \textit{UE})\} \end{array} \right\}$$

Codes

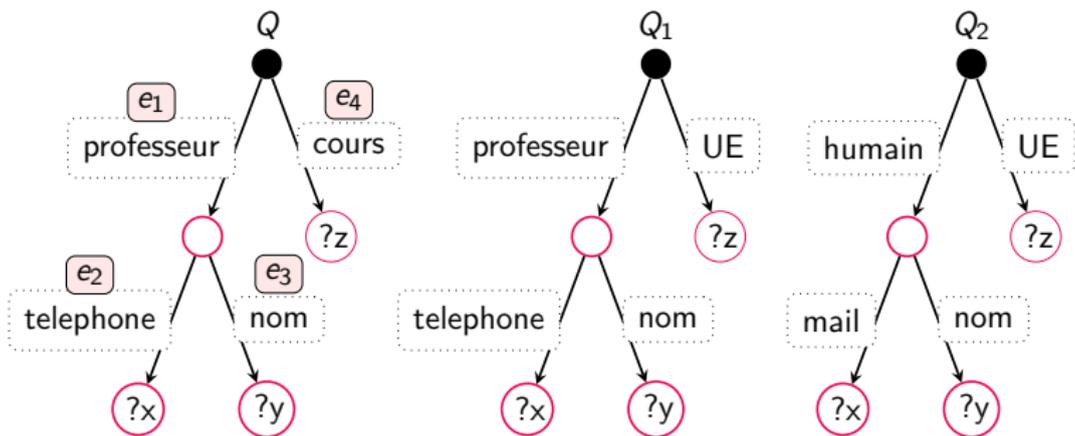
$c = (e_4, e_3, e_2, e_1)$

0000 = Q

1000 = Q_1

1012 = Q_2

$$C_c = \left\{ \begin{array}{l} CC_1 = \{(0, \textit{professeur}), (1, \textit{personne}), (2, \textit{humain})\} \\ CC_2 = \{(0, \textit{telephone}), (1, \textit{mail})\} \\ CC_3 = \{(0, \textit{nom})\} \\ CC_4 = \{(0, \textit{cours}), (1, \textit{UE})\} \end{array} \right\}$$



Équivalence code/entier

$$c = (e_4, e_3, e_2, e_1) \quad Cc = \left\{ \begin{array}{l} Cc_1 = \{(0, \textit{professeur}), (1, \textit{personne}), (2, \textit{humain})\} \\ Cc_2 = \{(0, \textit{telephone}), (1, \textit{mail})\} \\ Cc_3 = \{(0, \textit{nom})\} \\ Cc_4 = \{(0, \textit{cours}), (1, \textit{UE})\} \end{array} \right\}$$

$$c \equiv i$$

$$0000 \equiv 0$$

$$0001 \equiv 1$$

$$0002 \equiv 2$$

$$0010 \equiv 3$$

$$0011 \equiv 4$$

⋮

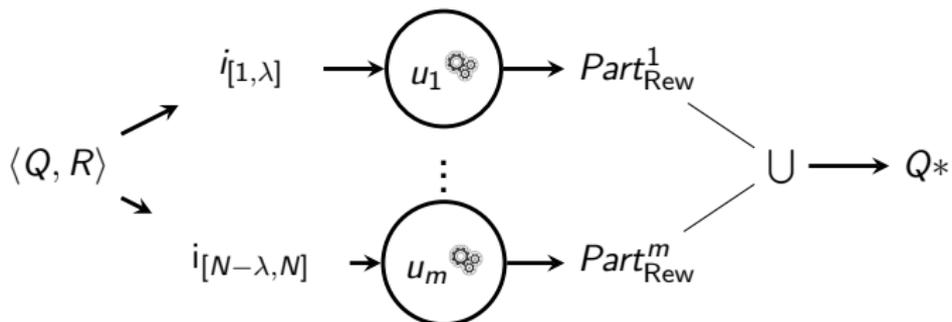
$$1012 \equiv 11$$

$$i = \begin{aligned} & e_4 \times (|Cc_1| \times |Cc_2| \times |Cc_3|) \\ & + e_3 \times (|Cc_1| \times |Cc_2|) \\ & + e_2 \times (|Cc_1|) \\ & + e_1 \end{aligned}$$

$$Q^* = [0, 11]$$

Parallélisation de la réécriture

- Idée : créer un partitionnement de l'espace des réécritures
- Unité de calcul associée avec un intervalle de taille $\lambda \approx \lceil N/m \rceil$



- Possibilité de paralléliser l'évaluation des réécritures sur les données.

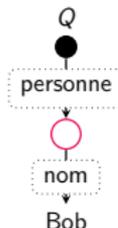
Traduction du langage de requête vers un langage cible

- Difficultés à assurer la correspondance entre le langage abstrait et un langage natif d'un key-value store (ex : mongodb)
- Non trivial : les key-value stores manquent de standardisation

Données

```
{  
  "personne" : {  
    "nom" : "Bob",  
    "tel" : "12-34-56-78-90"  
  }  
}
```

Requête



Traduction MongoDB "intuitive" ne sélectionne pas le record

```
db.find( { "personne" : { "nom" : "Bob" } } )
```

(*)Equality matches on the embedded document require an exact match of the specified document, including the field order.

État des travaux

- Traduction des requêtes arborescentes vers le langage de MongoDB
- Article démonstration accepté à BDA 2017 (comité de lecture)
« *Querying Key-Value Stores Under Simple Semantic Constraints : Rewriting and Parallelization* »
O. Rodriguez, C. Colomier, C. Rivière, R. Akbarinia, F. Ulliana
Démonstration présentée à Nancy en novembre 2017.

Premières expérimentations

Nombre de threads	1	2	3	4	5
Nombre de réécritures	Temps d'exécution (ms)				
18	10		1	1	1
36	1	2	1	1	1
108	15	3	3	2	2
324	6	4	3	3	3
648	6	5	4	3	3
1 296	9	9	7	11	7
2 592	25	15	14	18	9
18 432	180	186	144	152	34
93 312	3 662	1 234	883	525	590
124 416	1 407	1 055	1 516	678	930
186 624	12 144	2 475	1 490	2 363	1 103
209 952	1 370	826	1 293	1 045	523
248 832	7 400	5 452	3 577	1 141	1 107
746 496	32 478	14 393	13 268	17 206	13 660
839 808	58 435	30 302	24 687	20 361	20 004
1 259 712	34 084	22 343	13 602	19 425	11 567
5 971 968					
7 558 272					
22 674 816					
286 654 464					

Système d'intégration de données arborescentes

Threads 1 Collection Personnes_gen Synchronisation

Requête university_1 Règles university_1

```
{ "contactable" : { "$exists" : true } }
```

contact -> contactable

Calculer

Réécritures Requetage de la base Base de données Sortie commandes

```
{  
  "nom": "Luis",  
  "prenom": "Luparti",  
  "genre": "M",  
  "universite": {  
    "departement": "Services",  
    "nom": "Kursk State University"  
  }  
},
```

```
{  
  "nom": "Rosalynd",  
  "prenom": "York",  
  "genre": "F",  
  "universite": {  
    "departement": "Human Resources",  
    "nom": "University of Pitesti"  
  }  
},
```

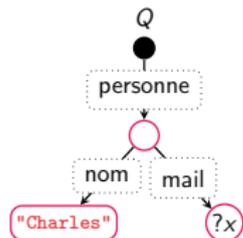
```
{  
  "nom": "Falkner",  
  "prenom": "Amthor",  
  "genre": "M",  
  "universite": {  
    "departement": "Engineering",  
    "nom": "GC University"  
  }  
},
```

```
{  
  "nom": "Aeriel",  
  "prenom": "McGlaughn",  
  "genre": "F",  
  "contactable": false,  
  "contact": [  
    {  
      "nom": "Nanine",  
      "prenom": "Ramel",  
      "genre": "F",  
      "universite": {  
        "departement": "Product Management"  
      }  
    },  
    {  
      "nom": "Margaretta",  
      "prenom": "Carreyette",  
      "genre": "F",  
      "contact": [  
        {  
          "nom": "Luis",  
          "prenom": "Luparti",  
          "genre": "M",  
          "universite": {  
            "departement": "Services",  
            "nom": "Kursk State University"  
          }  
        },  
        {  
          "nom": "Rosalynd",  
          "prenom": "York",  
          "genre": "F",  
          "universite": {  
            "departement": "Human Resources",  
            "nom": "University of Pitesti"  
          }  
        },  
        {  
          "nom": "Falkner",  
          "prenom": "Amthor",  
          "genre": "M",  
          "universite": {  
            "departement": "Engineering",  
            "nom": "GC University"  
          }  
        }  
      ]  
    }  
  ]  
}
```

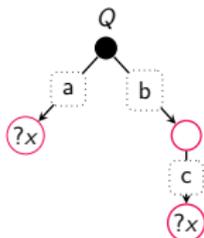
Objectifs à venir I

1) Vers un le langage d'interrogation plus expressif

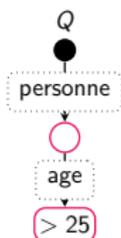
Actuellement



Jointure



Op. de comparaisons



2) Vers un langage de règles plus expressif

Actuellement

$(\forall) A \longrightarrow B$

$(\exists) A \longrightarrow \exists B$

Règles chemin ([MRU16, BBM⁺17])

$A_1 \dots A_n \longrightarrow [\exists] B_1 \dots B_m$

Point clé : la réécriture doit produire un ensemble fini de requêtes qui peuvent être évaluées dans le langage des sources

Objectifs à venir II

3) Étudier de nouvelles possibilités de traductions (pour le moment mongodb) pour confirmer que le langage de requête est concrètement utilisable sur les key-value stores.

▶ couchdb, dynamodb

4) Réfléchir aux optimisations en connaissant des contraintes sur les données (schéma).

5) Étudier la problématique de la fédération des bases de données.

