

Intégration de données

Exercices dirigés - corrigé

Pour chacun des cas suivants, donnez la description du médiateur (mapping entre le modèle global et les sources) et montrez comment les requêtes précisées sont traitées par le médiateur.

1. Global as view

Sources de données :

- S1:** *FilmTitre* (fid, titre)
FilmDetails (fid, année, genre, réalisateur)
Acteur (aid, nom, pays)
Joue (aid, fid)
- S2:** *FilmInfo* (titre, année, genre)
- S3:** *Realisateur* (nom, titre)

Modèle global du médiateur GAV:

- M:** *Film* (titre, année, genre, réalisateur)
Acteur (nom, titre)

Requêtes :

- a) Les titres et années des films de James Cameron.
- b) Les titres des films de James Cameron.
- c) Les acteurs du film Avatar.
- d) Les acteurs dirigés par James Cameron.

Solution

Mapping: $M = V(S1, S2, \dots)$

M: *Film* (t, a, g, r) :- **S1:** *FilmTitre* (fid, t) **S1:** *FilmDetails* (fid, a, g, r) U
S2: *FilmInfo* (t, a, g) **S3:** *Realisateur* (r, t)

M: *Acteur* (n, t) :- **S1:** *Acteur* (aid, n, p?) **S1:** *Joue* (aid, fid) **S1:** *FilmTitre* (fid, t)

Remarque: les mappings partiels, où pour un élément de M on ne peut retrouver qu'une partie des composantes, sont aussi possibles si on n'a pas de mapping complet pour l'élément. Ces mappings peuvent apporter des réponses aux requêtes qui ne concernent pas les composantes manquantes.

Requêtes

a) Les titres et années des films de James Cameron.

$Q(t, a) = M: \text{Film}(t, a, g', \text{"James Cameron"}) =$

S1: *FilmTitre* (fid, t) **S1:** *FilmDetails* (fid, a, g', "James Cameron") U

S2: *FilmInfo* (t, a, g') **S3:** *Realisateur* ("James Cameron", t)

$Q_1(t, a) = S1: FilmTitre (\underline{fid}, t) S1: FilmDetails (\underline{fid}, a, g', "James\ Cameron")$

$Q_2(t, a) = S2: FilmInfo (t, a, g')$

$Q_3(t) = S3: Realisateur ("James\ Cameron", t)$

En SQL :

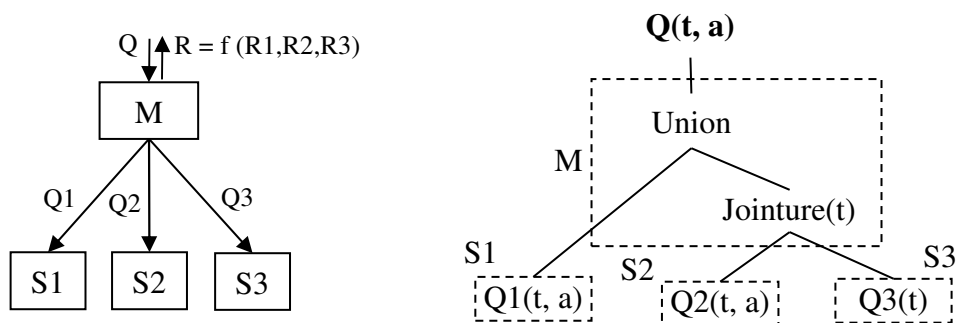
$Q(t, a) = \mathbf{select\ f.titre, f.ann\ee\ from\ M:Film\ f\ where\ f.r\realisateur = "James\ Cameron"}$

$Q_1(t, a) = \mathbf{select\ ft.titre, fd.ann\ee\ from\ S1:FilmTitre\ ft, S1:FilmDetails\ fd\ where\ fd.r\realisateur = 'James\ Cameron'\ and\ ft.fid=fd.fid}$

$Q_2(t, a) = \mathbf{select\ fi.titre, fi.ann\ee\ from\ S2:FilmInfo\ fi}$

$Q_3(t) = \mathbf{select\ r.titre\ from\ S3:Realisateur\ r\ where\ r.nom = 'James\ Cameron'}$

Plan d'exécution :



Optimisation : $Q_2(t, a)$ retourne tout S2, quand seulement les films de James Cameron sont utiles. Cela peut faire transiter sur le r\eseau une grande masse de donn\ees inutilement.

Am\elioration : r\realiser la jointure entre S2 et S3 directement dans S2. On r\ecup\ere d'abord les r\esultats de $Q_3(t)$ sous la forme d'une liste de titres de films LT et on envoie \a S2 la requ\ete modifi\ee suivante (en SQL) :

$Q_2'(t, a) = \mathbf{select\ fi.titre, fi.ann\ee\ from\ S2:FilmInfo\ fi\ where\ fi.titre\ in\ (LT)}$

b) Les titres des films de James Cameron.

$Q(t) = M: Film(t, a', g', "James\ Cameron") =$

$S1: FilmTitre (\underline{fid}, t) S1: FilmDetails (\underline{fid}, a', g', "James\ Cameron") \cup$

$S2: FilmInfo (t, a', g') S3: Realisateur ("James\ Cameron", t)$

$S2: FilmInfo (t, a', g')$ est ici inutile

$Q(t) = Q_1(t) \cup Q_3(t)$

c) Les acteurs du film Avatar.

$Q(n) = M: Acteur(n, "Avatar") =$

$S1: Acteur (\underline{aid}, n, p?) S1: Joue (\underline{aid}, \underline{fid}) S1: FilmTitre (\underline{fid}, "Avatar")$

$Q(n) = Q_1(n)$

d) Les acteurs dirig\es par James Cameron.

$Q(n) = M: Acteur(n, t) M: Film(t, a', g', "James\ Cameron") =$

$S1: Acteur (\underline{aid}, n, p?) S1: Joue (\underline{aid}, \underline{fid}) S1: FilmTitre (\underline{fid}, t) S1: FilmTitre (\underline{fid1}, t)$

$S1: FilmDetails (\underline{fid1}, a', g', "James\ Cameron") \cup$

$S1: Acteur (\underline{aid}, n, p?) S1: Joue (\underline{aid}, \underline{fid}) S1: FilmTitre (\underline{fid}, t) S2: FilmInfo (t, a', g')$

$S3: Realisateur ("James\ Cameron", t)$

Pour la première combinaison, on veut à travers la jointure sur le titre, avoir le même film. Mais ici, on dispose d'un moyen plus fort de s'en assurer (le titre n'étant pas une vraie clé pour les films) : le fid. Donc on peut mettre $fid=fid1$, pour avoir le même film.

La première combinaison devient :

S1: *Acteur* (aid, n, p?) **S1:** *Joue* (aid, fid) ~~**S1:** *FilmTitre* (fid, t)~~ **S1:** *FilmDetails* (fid, a', g', "James Cameron")

La présence de *FilmTitre* devient inutile, car le titre ne sert plus à faire le lien avec *Acteur*, le *fid* s'en charge.

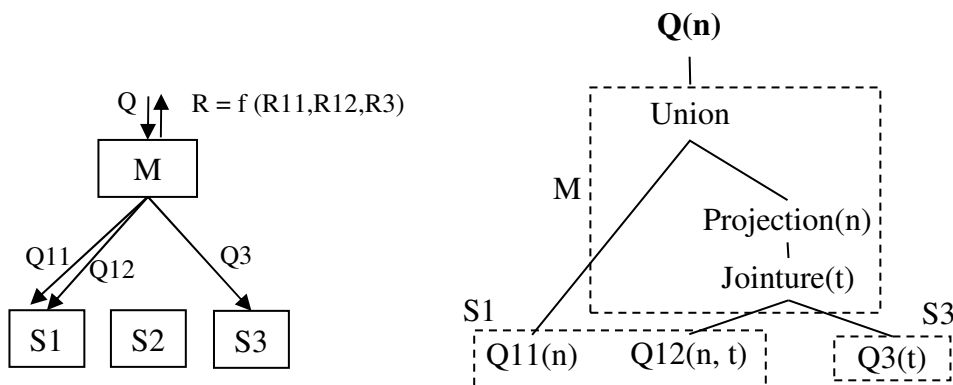
Résultat :

$Q(n) = M:Acteur(n, t) M:Film(t, a', g', "James Cameron") =$

S1: *Acteur* (aid, n, p?) **S1:** *Joue* (aid, fid) **S1:** *FilmDetails* (fid, a', g', "James Cameron") U

S1: *Acteur* (aid, n, p?) **S1:** *Joue* (aid, fid) **S1:** *FilmTitre* (fid, t) **S3:** *Realisateur* ("James Cameron", t)

$Q(n) = Q11(n) U Proj(n)(Q12(n, t) Q3(t))$



Remarque : le résultat est une union de combinaisons de sources (jointures entre sources).

Question : est-il possible d'avoir certaines de ces combinaisons redondantes, çàd leurs résultats sont déjà fournis par d'autres combinaisons ?

Réponse : sur le principe oui, une situation type est une jointure de type $S_i S_j$, comparée avec S_i tout seul. La jointure de S_i avec S_j apporte une contrainte supplémentaire aux données de S_i , donc la source S_i toute seule peut en principe fournir toute seule les résultats produits par la jointure $S_i S_j$.

Dans notre cas, on a une combinaison avec S_1 toute seule et une autre avec S_1 jointe à S_3 .

Dans S_1 tout seul, les films de James Cameron sont donnés par *FilmDetails*, dont on a dit qu'elle ne fournit d'information que pour une partie des films. Il est donc possible que certains des films de *FilmTitre* soient de James Cameron, sans qu'on l'indique dans *FilmDetails*. Donc la combinaison $S_1 S_3$ peut apporter de nouveaux résultats dans ce cas de figure, ce qui signifie qu'elle n'est pas redondante.

2. Local as view

Modèle global du médiateur LAV:

M: *Film* (titre, année, genre, réalisateur)

Acteur (nom, titre)

Sources de données :

S1: titre, années et réalisateurs de comédies

S2: acteurs avec le titre et année des films où ils ont joué

S3: réalisateurs, titres des films et acteurs qui y jouent

S4: titre et genre de films français d'avant 1970

Requêtes :

- Les titres et années des films de James Cameron.
- Les titres des films de James Cameron.
- Les acteurs du film Avatar.
- Les acteurs dirigés par James Cameron après 1980.

Solution

Mapping: $S_i \subseteq V_i(M)$

$S1(t, a, r) \subseteq M:Film(t, a, "comédie", r)$

$S2(n, t, a) \subseteq M:Acteur(n, t) M:Film(t, a, g', r')$

$S3(r, t, n) \subseteq M:Film(t, a', g', r) M:Acteur(n, t)$

$S4(t, g) \subseteq M:Film(t, a, g, r'), a < 1970$

Requêtes

- Les titres et années des films de James Cameron.

$Q(t, a) = M: Film(t, a, g', "James Cameron")$

S1 est la seule source qui permet de retourner le titre et l'année du film, tout en filtrant par le nom du réalisateur.

$Q(t, a) = S1(t, a, "James Cameron") = Q1(t, a)$

En SQL : $Q(t, a) = \text{select } S1.titre, S1.annee \text{ from } S1 \text{ where } S1.realisateur = 'James Cameron'$

- Les titres des films de James Cameron.

$Q(t) = M: Film(t, a', g', "James Cameron")$

S1 et S3 peuvent répondre à la requête.

$Q(t) = S1(t, a', "James Cameron") \cup S3("James Cameron", t, n') = Q1(t) \cup Q3(t)$

- Les acteurs du film Avatar.

$Q(n) = M:Acteur(n, "Avatar") = S2(n, "Avatar", a') \cup S3(r', "Avatar", n) = Q2(n) \cup Q3(n)$

d) Les acteurs dirigés par James Cameron.

$Q(n) = M:Acteur(n, t) M:Film(t, a', g', \text{"James Cameron"})$

On applique l'algorithme Bucket et on regarde pour chaque part de la requête quelles sources peuvent la fournir :

- $M:Acteur(n, t) : S2, S3$
- $M:Film(t, a', g', \text{"James Cameron"}) : S1, S3$

On génère toutes les combinaisons de sources possibles pour composer la requête :

- $S2(n, t, a') S1(t, a'', \text{"James Cameron"})$
- $S2(n, t, a') S3(\text{"James Cameron"}, t, n')$
- $S3(r', t, n) S1(t, a'', \text{"James Cameron"})$
- $S3(r', t, n) S3(\text{"James Cameron"}, t, n') = S3(\text{"James Cameron"}, t', n)$

La source S3 peut fournir des réponses toute seule, il faut voir si S2S3 ou S3S1 ne sont pas redondantes :

- $S2(n, t, a') S3(\text{"James Cameron"}, t, n')$ peut apporter des noms d'acteurs supplémentaires par rapport à ceux de S3 pour les films de James Cameron.
- $S3(r', t, n) S1(t, a'', \text{"James Cameron"})$: S1 apporte des titres de films de James Cameron, mais qui sont déjà dans S3 (jointure sur t), donc cette combinaison est redondante.

Résultat : $Q(n) = S3(\text{"James Cameron"}, t', n) \cup$
 $S2(n, t, a') S1(t, a'', \text{"James Cameron"}) \cup$
 $S2(n, t, a') S3(\text{"James Cameron"}, t, n')$

Pour effectuer la requête sur S2 une seule fois, on peut réécrire l'union de jointures en jointure d'unions :

$Q(n) = Q3(n) \cup \text{Projection}(n) (Q2(n, t) (Q1(t) \cup Q31(t)))$

