

Analyse théorique des algorithmes de jointure:

I. Jointure par Produit Cartésien:

Fonctionnement:

On suppose que l'on a deux relations: $R(X,Y)$, $S(Y,X)$, et que l'on applique la jointure suivante:
 $T = \text{Jointure}(R, S, R.Y = S.Y)$

1-lire la page i de R

2-lire la page j de S

3-effectuer les 4 matches possibles: $R.P_i.Y_0 \neq S.P_j.Y_0$

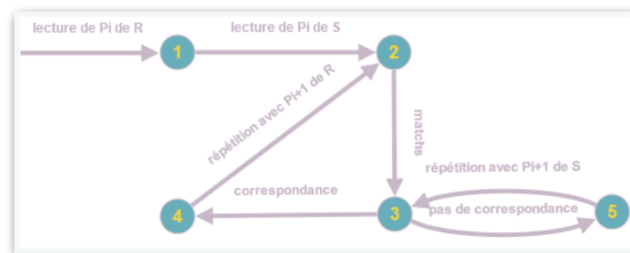
$$R.P_i.Y_1 \neq S.P_j.Y_0$$

$$R.P_i.Y_1 \neq S.P_j.Y_1$$

$$R.P_i.Y_0 \neq S.P_j.Y_1$$

4-s'il y a une correspondance alors l'écrire dans le résultat de la jointure (par exemple $T.append(R.P_i.X_0, R.P_i.Y_0, S.P_i.Y_0, R.P_i.X_0)$) et répéter les étapes 2,3,4,5 avec P_{i+1} de R

5-si pas de correspondance alors lire P_{j+1} de S et répéter les étapes 3,4,5



Graphe représentant le fonctionnement

de la jointure par Produit Cartésien

Avantage:

Convient pour des relations non triées et avec une distribution de données quelconque.

Inconvénient :

La durée d'exécution est généralement élevée, car elle nécessite de parcourir toutes les combinaisons possibles de tuples des relations.

Nombre de lectures et écritures disques:

-taille(R) lectures pour R

-taille(R)*taille(S) lectures pour S

-min(taille(R),taille(S)) écritures pour T

Total=taille(R)+taille(R)*taille(S)+min(taille(R),taille(S)) lectures/écritures disque.

II. Jointure par Tri-Fusion:

On suppose que l'on a deux relations déjà triées (et sans doublon) par ordre croissant en fonction de Y : $R'(X,Y)$, $S'(Y,X)$, et que l'on applique la même jointure: $T = \text{Jointure}(R', S', R'.Y = S'.Y)$

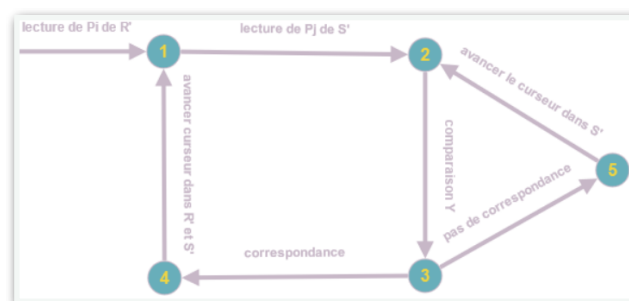
1-lire la page i de R'

2-lire la page j de S'

3- comparer $R'.P_i.Y_0$ avec $S'.P_j.Y_0$ et $R'.P_i.Y_1$ avec $S'.P_j.Y_1$

4-s'il y a une correspondance alors l'écrire dans le résultat de la jointure (T), avancer le curseur dans R' et dans S' et répéter les étapes 2,3,4,5.

5-s'il n'y a pas de correspondance, si $R'.P_i.Y_0 > S'.P_j.Y_0$ alors avancer le curseur dans S' ($Y(1)$) ou passer à P_{j+1} si déjà dans $Y(1)$) et répéter toutes les étapes.



Graphe représentant le fonctionnement
de la jointure par Tri-Fusion

Avantage:

Pas besoin de construire des index ou des tables de hachage.

Inconvénient :

Le tri initial peut être coûteux en termes de temps et de ressources dans le cas où il y a de grandes relations.

Nombre de lectures et écritures disques:

-taille(R) lectures pour R

-taille(S) lectures pour S

-min(taille(R),taille(S)) écritures pour T

Total=taille(R)+taille(S)+min(taille(R),taille(S)) lectures/écritures disque.

III. Jointure par Hachage Simple:

Fonctionnement:

On suppose que l'on reprend les deux relations: $R(X,Y)$, $S(Y,X)$ utilisées pour le produit cartésien (par conséquent non triées). on suppose également que l'on a une table de hachage "hachage" avant autant d'entrée que de résultat possible de la fonction de hachage, et que l'on applique la même jointure sur Y :
 $T = \text{Jointure}(R,S,R.Y=S.Y)$

Phase 1: Build

On suppose que l'on a choisit la fonction de hachage modulo 3.

1-lire P_i de R

2-pour chaque coupe (X,Y) calculer le modulo 3 de Y et l'écrire dans l'indice de la table de hachage correspondant ($\text{hachage}[Y\%3].\text{append}(R.P_i.X0, R.P_i.Y0, S.P_i.Y0, R.P_i.X0)$) en formant des blocs de deux éléments, répéter ces étapes avec $R.P_{i+1}$

Le résultat obtenu est une nouvelle version de la relation R où les Y ont été regroupés par la fonction de hachage.

Phase 2: Probe

1-lire la page i de S

2-calculer $S.P_j.Y0\%3$

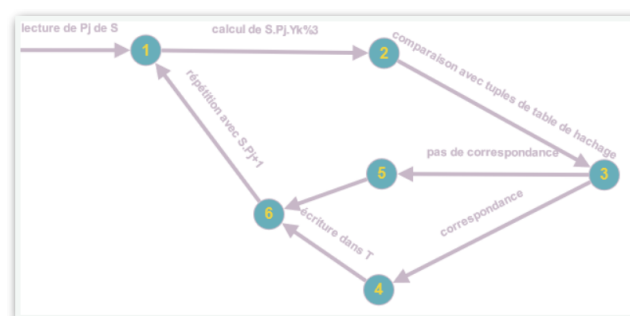
3-chercher un tuple pour lequel Y vaut $S.P_j.Y0\%3$ uniquement dans l'entrée $S.P_j.Y0\%3$ de "hachage"

4-s'il y a une correspondance l'écrire dans T entrées

5-si pas de correspondance alors arrêter la recherche sans vérifier les autres entrées

6-reproduire les memes étapes avec $S.P_{j+1}$ puis avec $S.P_{j+1}$

Graphe représentant le fonctionnement



de la jointure par Hachage Simple

Avantages:

Utilise une table de hachage pour une recherche rapide des correspondances.
Peut gérer efficacement les données non triées.

Inconvénient :

Nécessite l'élaboration d'une table de hachage, ce qui peut être coûteux en termes de mémoire et de ressources lors de la phase de « Build ».

Nombre de lectures et écritures disques:

-taille(R) lectures pour R

-taille(S) lectures pour S

-min(taille(R),taille(S)) écritures pour T

Total=taille(R)+taille(S)+min(taille(R),taille(S)) lectures/écritures disque.

IV. Jointure par Produit Cartésien Indexé:

Fonctionnement:

On suppose que l'on a deux relations: $R(X,Y)$, $S(Y,X)$, et que l'on applique la jointure suivante:
 $T = \text{Jointure}(R, S, R.Y = S.Y)$

Il faut mettre en oeuvre un index I sur l'attribut Y de S. Cette structure stocke les valeurs de Y triées dans l'ordre croissant et permet d'accéder à chacune de manière efficace et rapide.

Pour chaque tuple dans R, il faut chercher les correspondances dans l'index et combine les tuples correspondants de R et S pour former le tuple résultant.

1-lire la page i de R

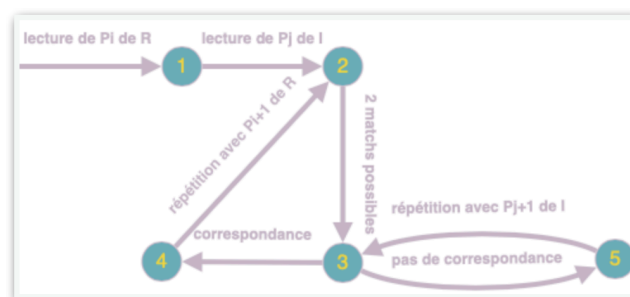
2-lire la valeur j de I

3-effectuer les 2 matches possibles: $R.P_i.Y_0 \neq I.P_j$

$R.P_i.Y_1 \neq I.P_j$

4-s'il y a une correspondance alors l'écrire dans le résultat de la jointure ($T.append(R.P_i.X_0, R.P_i.Y_0, S.(I.P_j).Y_0, S.(I.P_j).X_0)$ par exemple) et répéter les étapes 2,3,4,5 avec P_{i+1} de R

5-si pas de correspondance alors lire P_{j+1} de I et répéter les étapes 3,4,5



Graphe représentant le fonctionnement

de la jointure par Produit Cartésien Indexé

Avantages:

Utilise un index pour une recherche rapide des correspondances.

Peut gérer efficacement des relations non triées.

Inconvénient :

Nécessite l'élaboration d'un index, ce qui peut être coûteux en termes de mémoire et de ressources lors de la phase de prétraitement.

Nombre de lectures et écritures disques:

-taille(R) lectures pour R

-taille(S) lectures pour S

-min(taille(R),taille(S)) écritures pour T

Total=taille(R)+taille(S)+min(taille(R),taille(S)) lectures/écritures disque.

V. Jointure par Hachage Hybride:

....

IV. Comparaison:

Après avoir effectué une analyse théorique des différents algorithmes de jointures nous pouvons en conclure qu'en terme de lectures/écritures disques ce sont les algorithmes Tri-Fusion, Hachage simple et Produit Cartésien Indexé qui sont les plus efficaces.

En effet, leur exécution engendre le même nombre de lectures/écritures alors que le Produit Cartésien en engendre bien plus.

Néanmoins, certaines spécificités de ces trois algorithmes peuvent engendrer une distanciation par rapport aux autres. Ainsi, le Tri-Fusion semble être plus coûteux et avoir un temps d'exécution plus long que les deux autres car il est nécessaire de trier les données.

Le Hachage Simple et le Produit Cartésien Indexé nécessitent chacun l'élaboration d'une table de hachage ou d'un index ce qui est également coûteux.

L'analyse expérimentale devrait permettre de séparer ces algorithmes en fonction des distributions des données.