

Compte-rendu Projet SGBD

Perrin Matthieu - Gallay Jules

Décembre 2019



Sommaire

1	Présentation	3
2	Description des classes	3
2.1	Classe « Enregistrement »	5
2.2	Classe « Table »	6
2.3	Classe « Index »	6
2.4	Classe « Bloc »	7
2.5	Classe « Segment »	7
2.6	Classe « Buffer »	8
2.7	Autres classes	9
3	Description des algorithmes	10
3.1	Sort-merge	10
3.2	Key-lookup	12
3.3	Produit Cartésien	13
3.4	Omiecinski	14
3.5	Hachage	16
3.5.1	Fonction partition	16
3.5.2	Fonction hash	16
4	Interface	18
5	Interprétation des résultats	19

1 Présentation

Ce projet a pour but de comparer les coûts d'exécution de différents algorithmes de jointure dans des environnements variés. Le coût sera exprimé par le nombre de blocs lus et/ou écrits pour l'exécution des différents algorithmes, de plus il s'agira du coût réel d'exécution sur des tables que l'on aura créées.

2 Description des classes

Pour simuler ces différents tests, nous aurons besoins de briques de base que l'on utilisera pour calculer le coût des algorithmes. Il va s'agir des classes « Table », « Bloc », « Buffer », « Enregistrement », « Index » et « Segment ».

Ces classes fourniront les opérations nécessaires à leur mise en œuvre, des méthodes lectures/écritures seront mises à disposition par exemple.

Ci-dessous se trouve le diagramme de classe sur lequel nous nous sommes appuyés pour réaliser les différentes briques.

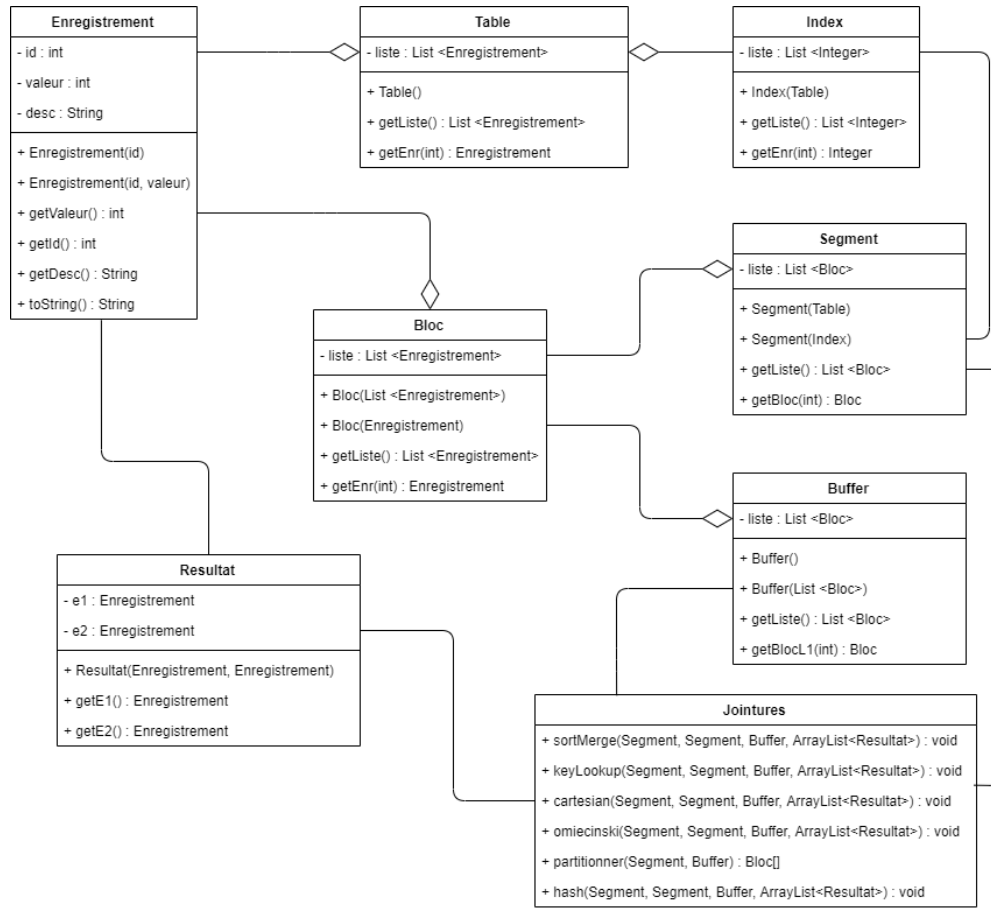


Diagramme de classe

2.1 Classe « Enregistrement »

La classe « Enregistrement » est une classe qui contiendra les valeurs que l'on comparera ensuite dans les algorithmes de jointure. Elle possède 3 attributs (id, desc et valeur) qui simule des informations que l'on peut inscrire dans un tuple. Les accesseurs ainsi que le constructeur n'ont rien de particulier, les valeurs sont définies par défaut aléatoirement et on crée un nouvel enregistrement en indiquant seulement l'id dans les paramètres. Un deuxième constructeur est disponible afin d'indiquer la valeur que l'on veut inscrire, et il est utilisé lors de la construction d'index sur l'attribut valeur.

```
1 //Constructeur
2 public Enregistrement(int p_id){
3     this.id=p_id;
4
5     this.valeur = 1 + (int)(Math.random() * 99);
6
7     byte[] array = new byte[7]; // length is bounded by 7
8     new Random().nextBytes(array);
9     this.desc = new String(array, Charset.forName("UTF-8"));
10 }
11
12 public Enregistrement(int p_id, int p_valeur){
13     this.id=p_id;
14     this.valeur = p_valeur;
15 }
```

Cette classe implémente l'interface Comparable et la méthode compareTo est surchargée sur l'attribut valeur. Cela permet de trier facilement les tables sur cet attribut.

```
1 @Override
2 public int compareTo(Enregistrement e) {
3     return this.valeur - e.getValeur();
4 }
```

Une méthode toString est aussi implémenté où l'on affiche les attributs de la classe, l'un au-dessus des autres.

```
1 @Override
2 public String toString(){
3     String res = ""+this.id+"\nValeur : "+this.valeur+"\nDescription : "+this.desc;
4     return res;
5 }
```

2.2 Classe « Table »

Cette classe possède un seul attribut qui est une ArrayList d'Enregistrement. Un accesseur permet de récupérer la liste entièrement et un autre permet de récupérer l'enregistrement selon sa position dans la liste.

```
1 //Attribut
2 private List <Enregistrement> liste = new ArrayList<Enregistrement
   >();
3
4 //Getter
5 public List<Enregistrement> getListe() {
6     return liste;
7 }
8
9 public Enregistrement getEnr(int i){
10    return this.liste.get(i);
11 }
```

Le constructeur est par défaut et créer une liste de 100 enregistrements qui auront des valeurs aléatoires.

```
1 //Constructeur
2 public Table(){
3     int i=0;
4     for(i=0; i<100; i++){
5         this.liste.add(new Enregistrement(i));
6     }
7 }
```

2.3 Classe « Index »

La classe « Index » possède une ArrayList d'Enregistrement comme attribut. Cette liste représente la liste des valeurs des enregistrements d'une table. Les accesseurs sont les mêmes que pour la classe « Table » mais pour la liste d'Integer et non la liste d'Enregistrement. Le constructeur récupère les valeurs de la table mise en paramètre et les inscrit dans la liste que l'on a pour attribut.

```
1 // Constructeur
2 public Index(Table t){
3     Iterator<Enregistrement> iter = t.getListe().iterator();
4     while (iter.hasNext()) {
5         Enregistrement element = iter.next();
6         int valeur = element.getValeur();
7         this.liste.add(valeur);
8     }
9 }
```

2.4 Classe « Bloc »

Cette classe est similaire à la classe « Table », la seule différence est le constructeur qui prend en paramètre une liste d'enregistrement. Dans ce projet, un bloc peut contenir au maximum 10 enregistrements.

```
1 // Constructeur
2 public Bloc(List <Enregistrement> p_liste){
3     this.liste=p_liste;
4 }
```

2.5 Classe « Segment »

La classe « Segment » possède un attribut ArrayList de Bloc avec des accesseurs pour récupérer la liste ainsi qu'un bloc selon sa position. Deux constructeurs sont mis à disposition, le premier avec une table en paramètre et le deuxième avec un index.

```
1 //Attribut
2 private List <Bloc> liste = new ArrayList<Bloc>();
3
4 //Getter
5 public List<Bloc> getListe() {
6     return liste;
7 }
8
9 public Bloc getBloc(int i){
10     return this.liste.get(i);
11 }
```

Ce constructeur permet de récupérer chaque enregistrement de la table pour ensuite créer un bloc avec une partie de cette liste d'enregistrement qui sera inséré dans le segment. Comme dit précédemment, un bloc ne contient que 10 enregistrements, donc il faut ajouter un nouveau bloc à la liste tous les 10 enregistrements

```
1 //Constructeur avec table
2 public Segment(Table t){
3     int count = t.getListe().size();
4     int i=0;
5     int j=0;
6     for (i=0; i<count; i=i+10){
7
8         List <Enregistrement> li_enr = new ArrayList<Enregistrement>
9         >();
10         for(j=0; j<10; j++){
11             li_enr.add(t.getListe().get(i+j));
12         }
13         Bloc b = new Bloc(li_enr);
14         this.liste.add(b);
15     }
```

```

14     }
15 }

```

Pour l'index c'est le même principe sauf lorsque l'on remplit la liste d'enregistrement pour créer le bloc. Comme l'index ne possède que la valeur il nous faut aussi un id que l'on va renseigner en déclarant un nouvel enregistrement. La description quant à elle n'est pas initialisée, pour simuler un index qui ne porte que sur le champ valeur.

```

1 //Constructeur avec index
2 public Segment(Index in){
3     int count = in.getListe().size();
4     int i=0;
5     int j=0;
6     for (i=0; i<count; i=i+10){
7
8         List <Enregistrement> li_enr = new ArrayList<Enregistrement>
9         >();
10        for(j=0; j<10; j++){
11            li_enr.add(new Enregistrement(i+j, in.getListe().get(i+
12            j)));
13        }
14        Bloc b = new Bloc(li_enr);
15        this.liste.add(b);
16    }
17 }

```

2.6 Classe « Buffer »

Cette dernière brique est le buffer qui va être utilisé pour stocker les blocs lus après la jointure que l'on aura choisi d'effectuer. On y retrouve donc une arrayList de Bloc avec ces accesseurs.

```

1 //Attribut
2 private List <Bloc> liste;
3
4 //Getter
5 public List<Bloc> getListe() {
6     return liste;
7 }
8
9 public Bloc getBlocL1(int i){
10     return this.liste.get(i);
11 }
12
13 //Constructeurs
14 public Buffer(List <Bloc> l1){
15     this.liste= l1;
16 }
17
18 public Buffer(){
19     this.liste= new ArrayList<Bloc>();
20 }

```


2.7 Autres classes

Notre projet comporte d'autres classes dont la première est la classe « Resultat » qui possède 2 attributs « Enregistrement » que l'on va utiliser afin de récupérer les enregistrements qui ont une valeur commune lors des jointures.

```
1 //Attributs
2 private Enregistrement e1;
3 private Enregistrement e2;
4
5 //Getter
6 public Enregistrement getE1() {
7     return e1;
8 }
9
10 public Enregistrement getE2() {
11     return e2;
12 }
13
14 //Constructeur
15 public Resultat(Enregistrement en1, Enregistrement en2){
16     this.e1=en1;
17     this.e2=en2;
18 }
```

La deuxième est la classe « Jointures » qui possède tous les algorithmes que nous allons voir après.

On va y retrouver les 5 algorithmes de jointures : Key-lookup, Sort-merge, Produit Cartésien, Hashage et Omiecinski

On a décidé d'effectuer le tri avant la fonction et non à l'intérieur de celle-ci pour plus de facilité. On sélectionnera si les listes seront triées ou non dans le menu

3 Description des algorithmes

3.1 Sort-merge

L'algorithme Sort-merge est composé en 2 étapes : trié les deux relations sur le critère de jointure puis fusionner ces relations pour enfin calculer le coût.

On a décidé d'effectuer le tri avant la fonction et non à l'intérieur de celle-ci pour plus de facilité. On sélectionnera si les listes seront triées ou non dans le menu. Nous trions les relations grâce à la méthode sort de la classe List, par conséquent nous ne pouvons pas calculer le coût du tri, car les relations sont triées avant d'être « convertie » en segments et en blocs.

Notre fonction possède 4 paramètres : 2 segments, 1 buffer et une liste résultat. Premièrement on récupère la liste de bloc pour chaque segment et on déclare 2 variables pos1 et pos2 qui vont donner la position dans les blocs afin de pouvoir tout comparer. Les variables stock1 et stock2 serviront de stockage de cette position, pour vérifier si on entre dans un nouveau bloc. Elles vont nous permettre de calculer le coût de l'algorithme

```
1 //Attributs
2 List <Bloc> lb1 = s1.getListe();
3 List <Bloc> lb2 = s2.getListe();
4
5 int pos1=0;
6 int pos2=0;
7
8 int stock1=0;
9 int stock2=0;
10
11 //Boucle
12 ...
```

Une boucle while est utilisé afin de pouvoir comparer chaque valeur des 2 segments. On récupère l'Enregistrement ainsi que la valeur de celui-ci pour les deux segments que l'on va ensuite comparer.

```
1 //Attributs
2 ...
3
4 //Boucle
5 while (pos1 <= s1.getListe().size()*s1.getBloc(0).getListe().size()
6         -1 && pos2 <= s2.getListe().size()*s2.getBloc(0).getListe().
7         size()-1){
8
9     Enregistrement enr1 = s1.getBloc(pos1/10).getEnr(pos1%10);
10    Enregistrement enr2 = s2.getBloc(pos2/10).getEnr(pos2%10);
11
12    int val1=s1.getBloc(pos1/10).getEnr(pos1%10).getValeur();
13    int val2=s2.getBloc(pos2/10).getEnr(pos2%10).getValeur();
14
15    //Test des valeurs
16    ...
}
```

Si les deux valeurs sont égales alors on ajoute le couple de valeurs dans la liste Résultat, on stocke la position que l'on incrémente ensuite et pour finir on ajoute le bloc au buffer.

```

1 //Boucle
2 ...
3
4 //Verification des resultat
5 if(val1 == val2){
6     Resultat r = new Resultat(enr1, enr2);
7     liste.add(r);
8     stock2 = pos2/10;
9     pos2++;
10    if(pos2/10>stock2){
11        b.getListe().add(s2.getBloc(stock2));
12    }
13 }
14
15 //Suite des verifications
16 ...

```

Sinon si la valeur du segment 1 est supérieur à celle du segment 2 alors on incrémente juste la position et on ajoute le bloc précédent au buffer si on change de bloc. Et inversement si la valeur de segment 2 est supérieur à celle du segment 1.

```

1 //Boucle
2 ...
3
4 //Suite des verifications
5 else if(val1 > val2){
6     stock2 = pos2/10;
7     pos2++;
8     if(pos2/10>stock2){
9         b.getListe().add(s2.getBloc(stock2));
10    };
11 }
12 else{
13     stock1 = pos1/10;
14     pos1++;
15     if(pos1/10>stock1){
16         b.getListe().add(s1.getBloc(stock1));
17     }
18 }

```

Enfin on ajoute les derniers blocs au buffer si la position finale est supérieure à celle stockée en mémoire

```

1 //En dehors de la boucle
2 if(pos1/10>stock1){
3     b.getListe().add(s1.getBloc((pos1-1)/10));
4 }
5 if(pos2/10>stock2){
6     b.getListe().add(s2.getBloc((pos2-1)/10));
7 }

```

3.2 Key-lookup

L'algorithme Key-lookup se compose de boucles imbriquées pour comparer les tuples et les paramètres sont les mêmes que pour le premier algorithme. La première boucle récupère les valeurs du segment 1 qu'on ajoute au buffer puis la deuxième boucle récupère les valeurs du 2e segment qui sont aussi ajoutées au buffer.

```
1 //1ere boucle
2 for(i=0; i<s1.getListe().size()*s1.getBloc(0).getListe().size(); i
  ++){
3   //recupere enregistrement
4   Enregistrement enr1 = s1.getBloc(i/10).getEnr(i%10);
5   int val1=s1.getBloc(i/10).getEnr(i%10).getValeur();
6   //Stoque dans le bloc
7   if(i/10>stock){
8     b.getListe().add(s1.getBloc(stock));
9   }
10
11 //2e boucle
12 for(j=0; j<s2.getListe().size()*s2.getBloc(0).getListe().size()
  ; j++){
13   //Stoque dans le bloc
14   if(j/10>stockj){
15     b.getListe().add(s2.getBloc(stockj));
16   }
17   //recupere enregistrement
18   Enregistrement enr2 = s2.getBloc(j/10).getEnr(j%10);
19   int val2=s2.getBloc(j/10).getEnr(j%10).getValeur();
20
21   //Verification des resultats
22   ...
23 }
24 ...
25 }
```

Ensuite on compare les valeurs et si elles sont égales alors on les ajoute à la liste Résultat. On ajoute les derniers blocs à la fin comme pour Sort-merge.

```
1 //Boucle
2 ...
3
4 //Verification des resultats
5 if(val1 == val2){
6   Resultat r = new Resultat(enr1, enr2);
7   liste.add(r);
8 }
```

3.3 Produit Cartésien

Le principe de cet algorithme est de charger plusieurs blocs de la première relation et un bloc de la deuxième pour ensuite le comparer avec les autres blocs. Cette étape est répétée plusieurs fois jusqu'à ce que tous les blocs soient lus.

Pour ce qui est de notre fonction, les paramètres sont toujours les mêmes avec plusieurs variables déclarées comme par exemple une variable pour connaître le nombre de blocs que l'on chargera en mémoire centrale pour la première relation (qu'on appelle "capa" dans notre code).

Une première boucle permet de répéter le nombre de fois nécessaire pour que tous blocs du segment 1 soit chargés et comparés.

```
1 //On repete size/capa fois (ici 2)
2 for(i=0; i<(s1.getListe().size())/capa; i++){
3
4     //on charge les blocs a lire de r1
5     lb1.clear();
6     for(j=0; j<capa; j++){
7         lb1.add(s1.getBloc(i*5+j));
8     }
9
10    //Boucle pour recuperer les autres blocs
11    ...
12
13 }
```

Une deuxième boucle charge tous les blocs du segment 2 un à un pour les comparer avec les blocs du segment 1 chargés. Lorsque deux valeurs sont égales le traitement est similaire que les autres algorithmes.

```
1 //Boucle
2 ...
3
4 //on parcourt tous les blocs de r2 pour les comparer avec les 5
   blocs
5 for(k=0; k<s2.getListe().size()*s2.getBloc(0).getListe().size(); k
   ++){
6     if(k/10>stock_k){
7         b.getListe().add(s2.getBloc(stock_k));
8     }
9
10    //Recuperation des valeurs puis comparaison
11    ...
12
13 }
```

3.4 Omiecinski

Cet algorithme ressemble au Sort-merge mais ici des indexes secondaires sont utilisés. Premièrement, on trie les indexes sur leurs attributs, ce qui est effectué avant de lancer la fonction, puis on effectue les comparaisons entre les valeurs.

La fonction possède deux boucles imbriquées qui permet de parcourir chaque bloc que l'on va ensuite traiter. La première boucle parcourt les éléments de la première relation.

```
1 //Boucle 1
2 while(cpt1 < s1.getListe().size()*s1.getBloc(0).getListe().size()
3     -1){
4     Enregistrement enr1 = s1.getBloc(cpt1/10).getEnr(cpt1%10);
5     Enregistrement enr2 = s2.getBloc(cpt2/10).getEnr(cpt2%10);
6     int val1=s1.getBloc(cpt1/10).getEnr(cpt1%10).getValeur();
7     int val2=s2.getBloc(cpt2/10).getEnr(cpt2%10).getValeur();
8
9     //Boucle 2
10    ...
11
12    //Envoie des valeurs a la liste Resultat
13    ...
14
15 }
```

La deuxième boucle parcourt non seulement les éléments de la relation 2 mais vérifie si les éléments sont différents.

```
1 //A l'interieur de la boucle 1
2 ...
3
4 //Boucle 2
5 while(val1 != val2 && cpt2<s2.getListe().size()*s2.getBloc(0).
6     getListe().size()-1){
7     stock2=cpt2/10;
8     cpt2++;
9     if(cpt2/10>stock2){
10        b.getListe().add(s2.getBloc(stock2));
11    }
12    val2=s2.getBloc(cpt2/10).getEnr(cpt2%10).getValeur();
13    enr2 = s2.getBloc(cpt2/10).getEnr(cpt2%10);
14
15 }
16
17 //Envoie des valeurs a la liste Resultat
18 ...
```

Si les valeurs sont égales alors on sort de la boucle pour les envoyer dans la liste Résultat.

```
1 //Boucle 2
2 ...
3
4 //Envoie des valeurs a la liste Resultat
5 Resultat r = new Resultat(enr1, enr2);
6 liste.add(r);
7
8 stock1=cpt1/10;
9 cpt1++;
10 if(cpt1/10>stock1){
11     b.getListe().add(s1.getBloc(stock1));
12 }
13 cpt2=0;
```

3.5 Hachage

3.5.1 Fonction partition

On crée d'abord un tableau de Blocs qui récupérera les enregistrements après partition puis on parcourt les éléments du segment mis en paramètre.

Chaque enregistrement est récupéré puis enregistré dans le tableau de Blocs, stocké à l'indice de leur valeur. Par exemple un enregistrement de valeur 47 sera dans le bloc à l'indice 47 dans le tableau. Une fois tout ceci effectué on renvoie le tableau de Blocs

```
1 Bloc tab[] = new Bloc[100];
2
3 //on stocke les enregistrements a l'indice de leur valeur
4 for(i=0; i<s1.getListe().size()*s1.getBloc(0).getListe().size(); i++){
5
6     if(i/10>stock){
7         b.getListe().add(s1.getBloc(stock));
8     }
9     Enregistrement enr1 = s1.getBloc(i/10).getEnr(i%10);
10    int val1=s1.getBloc(i/10).getEnr(i%10).getValeur();
11    tab[val1]= new Bloc(enr1);
12
13    stock=i/10;
14 }
15 if(i/10>stock){
16     b.getListe().add(s1.getBloc((i-1)/10));
17 }
18
19 return tab
```

3.5.2 Fonction hash

Dans cette fonction on appelle la fonction ci-dessus pour chaque segment en paramètre puis une première boucle récupère un élément du premier et du deuxième tableau où chaque bloc est ajouté au buffer.

```
1 Bloc[] tab1= partitionner(s1,b);
2 Bloc[] tab2= partitionner(s2,b);
3
4 //Boucle 1
5 for(i=0; i<100; i++){
6     Bloc b1 = tab1[i];
7     Bloc b2 = tab2[i];
8
9     b.getListe().add(b1);
10    b.getListe().add(b2);
11
12    //Recuperation de la taille du tableau b2 pour la deuxi me
    boucle
13    ...
14 //Boucle 2
15    ...
16 }
```


Une deuxième boucle va ajouter les enregistrements du 2e bloc à une variable puis une dernière boucle récupère les enregistrements du 1er bloc.

```
1 //Boucle 1
2 ...
3
4 //Boucle 2
5 for(j=0; j<sizeb2; j++){
6     Enregistrement enr2 = b2.getEnr(j);
7
8     //Recuperation de la taille du tableau b1 pour la derniere
9     //boucle
10    ...
11    //Boucle 3
12    ...
13 }
```

Ces enregistrements seront envoyés dans la liste Résultat car les enregistrements partitionnés précédemment sont stockés à l'indice de leur valeur c'est-à-dire que deux enregistrements au même indice sont forcément égaux.

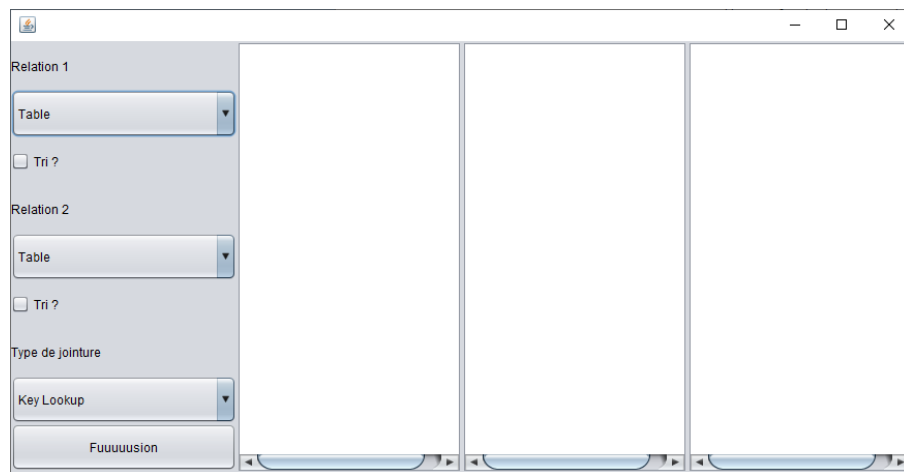
```
1 //Boucle 2
2 ...
3
4 //Boucle 3
5 for(k=0; k<sizeb1; k++){
6     Enregistrement enr1 = b1.getEnr(k);
7     Resultat r = new Resultat(enr1, enr2);
8     liste.add(r);
9 }
```

4 Interface

Pour pouvoir effectuer tous les tests nécessaires, nous avons décidé de faire une interface graphique avec NetBeans pour plus de facilité et de rapidité.

Cette interface se compose d'un menu à gauche où l'on peut sélectionner une table ou un index pour chaque relation. On peut aussi choisir si la relation sera triée ou non, sachant que la jointure Sort-merge et Omiecinski n'affichera aucun résultat si les deux relations ne sont pas triées.

Enfin on choisit le type de jointure que l'on veut parmi les 5 disponibles. 3 zones de textes afficheront les informations concernant les relations, la première pour les éléments de la relation 1, la deuxième pour les éléments de la relation 2 et enfin la dernière pour les résultats après jointure.



5 Interprétation des résultats

Key-Lookup : 1010 blocs lus

Ce coût très important s'explique par le fait que l'on ne puisse pas faire facilement une sélection sur la relation R2 avec l'index secondaire car nos index sont modélisés comme des tables. Pour une jointure de ce type, avec notre algorithme, le coût est de $Br1 + Tr1 * Br2$ soit $10 + 100 * 10 = 1010$.

Sort-merge : 20 blocs lus

Ce coût paraît très attractif, mais il n'est pas représentatif de la réalité car nous n'incluons pas le coût du tri des 2 relations dans l'algorithme. Le tri est géré en amont et nous n'avons donc ici que le coût de la jointure des relations triées qui est de $Br1 + Br2$ soit 20.

Produit Cartésien : 30 blocs

Ce coût est similaire au coût obtenu en condition réelle pour cet algorithme. Notre mémoire centrale accueille 5 blocs de R1 et les compare séquentiellement avec 1 blocs de R2, elle a donc une capacité de 6 blocs. On retrouve ce coût grâce à la formule du cours qui est $(Br1 / (M - 1)) * Br2 + Br1$ soit ici $(10 / 5) * 10 + 10 = 30$.

Hachage : 220 blocs lus

Ce coût peut paraître élevé, mais il est en fait cohérent avec les conditions réelles. Ce qui le fait autant monter est la fonction de partitionnement qui n'est pas optimisée car la liste de blocs qui en résulte a autant d'indices que de valeurs possibles pour l'attribut valeur d'un enregistrement.

Notre fonction de partitionnement a un coût de $Br1$, on l'utilise deux fois donc $Br1 + Br2$ soit 20 dans notre cas.

Ensuite la fonction de jointure parcourt les deux tableaux de blocs partitionnées, donc ces tableaux font 100 cases, et le coût est donc de 200. Notre coût total est donc bien de 220.

Omiecinski : 20 blocs lus

Tout comme l'algorithme Sort-merge, ici le coût paraît attractif car nous ne comptons pas le tri des relations. Le tri est géré en amont et nous n'avons donc ici que le coût de la jointure des relations triées qui est de $Br1 + Br2$ soit 20.

Ces résultats sont ceux observés avec nos choix de modélisation, ils ne reflètent donc pas forcément les coûts que nous pouvons observer sur un SGBD comme Oracle. Dans notre cas, l'algorithme du produit cartésien semble le plus attractif si l'on exclu les jointures sort merge et Omiecinski du fait du coût du tri absent.

A cause de nos choix de modélisation, l'algorithme key-lookup est à éviter absolument car il parcourt beaucoup trop de blocs.

Si nous parvenons à améliorer la fonction de partitionnement, l'algorithme de hachage peut devenir lui aussi très intéressant