

# Projet LSINF1121 - Algorithmique et structures de données

-

## Rapport intermédiaire Mission 4

Groupe 26

Laurian DETIFFE  
(6380-12-00)

Sundeeep DHILLON  
(6401-11-00)

Alexis MACQ  
(5910-12-00)

Xavier PÉRIGNON  
(8025-11-00)

Thibaut PIQUARD  
(4634-13-00)

Thomas WYCKMANS  
(3601-12-00)



Année académique 2015-2016

## Questions et réponses

1. Citez au moins quatre implémentations différentes d'un dictionnaire (table de symboles). Précisez, dans chaque cas, quelles sont les propriétés principales de ces implémentations. Dans quel(s) cas s'avèrent-elles intéressantes ? Quelles sont les complexités calculatoires de leurs principales méthodes ? (Xavier)

(a) **Recherche séquentielle (liste désordonnée) :**

Une option simple pour la structure de données d'une table de symboles est une liste chaînée de noeuds qui contiennent des clés et des valeurs. Le principe de l'algorithme est qu'il parcourt la liste en comparant la clé de recherche avec la clé de chaque noeud dans la liste. Si les clés concordent, il retourne la valeur associée. Sinon, il retourne *null*. La mise en œuvre d'une liste liée à la recherche séquentielle est trop lente pour qu'elle puisse être utilisée pour résoudre d'énormes problèmes. Cependant, elle est idéale pour les petits problèmes.

(b) **Binary tree search (BST) :**

Un arbre de recherche binaire est un arbre binaire où chaque noeud a une clé (et une valeur associée) et qui satisfait la restriction que la clé dans un noeud quelconque est plus grande que celle de tous les noeuds du sous-arbre gauche de ce noeud et plus petite que celle de tous les noeuds du sous-arbre droit de ce noeud. La mise en œuvre de cette structure est très facile à implémenter. Cependant, si l'arbre ressemble à une liste chaînée, cela risque d'augmenter fortement la complexité.

(c) **Separate chaining (tableau de listes) :**

Dans la méthode dite de chaînage séparé, chaque bucket est indépendant, et a une sorte de liste des entrées avec le même index. Le temps d'exécution d'une opération d'une table de hachage est le temps de trouver le bucket (qui est constant) plus le temps de l'opération concernant la liste. Dans une bonne table de hachage, chaque bucket a 0 ou 1 entrée, parfois 2 ou 3, mais rarement plus que cela. Par conséquent, les structures sont efficaces dans le temps et l'espace. Cependant, ce procédé hérite également des inconvénients des listes liées. En effet, lors du stockage de petites clés et de valeurs, la surcharge de l'espace du prochain pointeur dans chaque entrée peut être importante.

(d) **Linear probing (tableaux parallèles) :**

Linear probing est réalisée en utilisant deux valeurs : une étant comme une valeur de départ et une autre étant comme un intervalle entre les valeurs successives en arithmétique modulaire. La seconde valeur, qui est la même pour toutes les clés et connue sous le nom *stepsize*, est ajoutée à plusieurs reprises à la valeur de départ jusqu'à ce qu'un espace libre est trouvé, ou que toute la table soit parcourue.

$$newLocation = (startingValue + stepSize) \% arraySize$$

Cet algorithme offre une bonne mise en cache de la mémoire (si *stepsize* est égal à 1), grâce à la bonne localité de référence.

## Résumé :

algorithm (data structure)	worst-case cost (after N inserts)		average-case cost (after N random inserts)		key interface	memory (bytes)
	search	insert	search hit	insert		
<i>sequential search</i> (unordered list)	$N$	$N$	$N/2$	$N$	<code>equals()</code>	$48 N$
<i>binary tree search</i> (BST)	$N$	$N$	$1.39 \lg N$	$1.39 \lg N$	<code>compareTo()</code>	$64 N$
<i>separate chaining</i> (array of lists)	$< \lg N$	$< \lg N$	$N/(2M)$	$N/M$	<code>equals()</code> <code>hashCode()</code>	$48 N + 64 M$
<i>linear probing</i> (parallel arrays)	$c \lg N$	$c \lg N$	$< 1.50$	$< 2.50$	<code>equals()</code> <code>hashCode()</code>	between $32 N$ and $128 N$

2. Souvenez vous de la question suivante proposée en bilan sur la mission sur les tris : Étant donné un ensemble  $S$  de taille  $n$ , et un nombre  $x$ . Décrivez un algorithme efficace utilisant une HashTable pour trouver s'il existe une paire  $(a, b)$  avec  $a \in S$ ,  $b \in S$  telle que  $a + b = x$ . Quelle est la complexité de votre algorithme ? Est-elle meilleure que votre solution qui utilisait un tri ? (Xavier)

Concernant le bilan sur la mission sur les tris, la solution était de faire un tri efficace de l'ensemble  $S$  (par exemple Quicksort) et ensuite de tester l'addition du premier élément  $a$  avec le dernier élément  $b$ . Si le résultat de cette addition était inférieur à  $x$ , alors on testait l'addition de  $(a + 1)$  avec  $b$ , et ainsi de suite. Sinon si le résultat de l'addition était supérieur à  $x$ , alors on testait l'addition de  $a$  avec  $(b - 1)$ , et ainsi de suite.

Pour résoudre le même problème avec une HashTable, on pourrait penser à insérer tous les éléments de l'ensemble  $S$  dans une HashSet. Cette classe implémente l'interface Set, en utilisant une HashTable. HashSet est implémenté comme une HashMap, dont les clés sont les éléments du HashSet et les valeurs sont toutes une même valeur présente. Cette implémentation offre des performances constantes pour les opérations `add(T t)`, `remove(T t)`, `contains(T t)` et `size()`. Il suffira donc de vérifier si pour chaque élément  $a$  de l'ensemble  $S$ , il existe  $b$  dans la HashSet tel que  $b = x - a$ . La complexité de cet algorithme est donc de  $O(1) + O(n)$  (parcours des éléments de  $S$ ). Ce qui est plus rapide que l'algorithme de la mission sur les tris.

3. Démontrez que  $(a + b) \% M$  est équivalent à  $((a \% M) + b) \% M$ . En quoi cette propriété peut être utile pour construire une fonction de hachage sur les String. Expliquez comment Java calcule une fonction de hachage sur les String ? Quelle est

la complexité pour calculer 1 fois et N fois le hashcode d'un String. (Xavier)

Prenons  $G = (a + b) \% M$  et  $D = ((a \% M) + b) \% M$ .

On doit démontrer que  $G = D$ .

Nous pouvons écrire :

$a = M * Q1 + R1$  where  $0 \leq R1 < M$  où  $Q1$  est un entier.  $a \bmod M = R1$

$b = M * Q2 + R2$  where  $0 \leq R2 < M$  où  $Q2$  est un entier.  $b \bmod M = R2$

$$(a + b) = M * (Q1 + Q2) + R1 + R2$$

$$G = (a + b) \bmod M$$

$$G = (M * (Q1 + Q2) + R1 + R2) \bmod M$$

On peut supprimer les multiples de M lorsqu'on prend mod M,

$$G = (R1 + R2) \bmod M$$

$$D = (a \bmod M + b) \bmod M$$

$$D = (R1 + R2) \bmod M$$

Et donc,

$$G = D = (R1 + R2) \bmod M$$

Cette propriété est utile pour construire une fonction de hashage sur les String car il n'est pas nécessaire de calculer le modulo de chaque valeur des caractères du String. En effet, il suffira de calculer la somme de la valeur des caractères et utiliser une seule fois le modulo à la fin.

La fonction de hashage pour un String est implémentée comme suit :

$$s[0] * 31^{(n-1)} + s[1] * 31^{(n-2)} + \dots + s[n-1]$$

Dans cette fonction, le String  $s$  a  $n$  caractères. Elle calcule la somme des valeurs de chaque caractère, chacun multiplié par  $31^r$  avec  $r$  étant l'emplacement inversé du caractère dans le String. La complexité du calcul de hashcode d'un String est de  $O(n)$  puisqu'il faut parcourir tous les caractères du String. Mais une fois ce calcul fait, il n'est plus nécessaire de recalculer le hashcode de ce String.

4.

5. Java fournit la classe `java.util.Hashtable` comme implémentation de l'interface `java.util.Map`. Pouvez-vous déterminer précisément de quelle variante de table de hachage il s'agit ? Java fournit-il d'autres implémentations de l'interface `Map` ? Faites un diagramme qui représente les interfaces et les classes qui se rapportent à `Map` et précisez ce qui, dans chaque cas, le caractérise. Qu'est-ce qui peut servir de clef pour une `Hashtable` en Java ? Soyez précis. (Sundeeep)

Dans la classe `java.util.Hashtable` qui implémente l'interface `java.util.Map`, utilise la méthode du separate chaining pour la gestion des collisions dans la table

de hachage. Plusieurs constructeurs sont mis à disposition. Le premier, ne prenant aucun paramètre et qui construit une nouvelle table de hachage avec une capacité initiale et un facteur de charge par défaut, qui sont respectivement 11 et 0.75. Le second constructeur permet de choisir la capacité initiale et d'utiliser le facteur de charge par défaut. Le troisième, quant à lui, permet de choisir la capacité initiale et le facteur de charge. Cette version de la table de hachage proposée par `java.util.Hashtable` est synchronisée. Dans une `HashTable`, n'importe quel objet non-null peut être utilisé comme clé (ou valeur).

Le diagramme de classe que nous avons dessiné se trouve en annexe, voir figure ???. Nous avons délibérément ignoré les classes qui fournissaient des fonctionnalités qui n'avaient plus vraiment de lien avec un dictionnaire, et la classe `AbstractMap` qui aurait surchargé le schéma de flèches sans vraiment apporter d'information additionnelle.

6. Qu'entend-on par la notion de collision dans une table de hachage ? Les collisions ont-elles une influence sur la complexité des opérations ? Si oui, quelle(s) opération(s) avec quelle(s) complexité(s), sinon précisez pourquoi. (Thomas)

Le fait de créer une valeur de hachage à partir d'une clé peut engendrer un problème de « collision », c'est-à-dire que deux clés différentes, voire davantage, pourront se retrouver associées à la même valeur de hachage et donc à la même case dans le « tableau » (la fonction n'est pas injective). Pour diminuer les risques de collisions, il faut donc premièrement choisir avec soin sa fonction de hachage. Ensuite, un mécanisme de résolution des collisions sera à implémenter.

Tout comme les tableaux ordinaires, les tables de hachage permettent un accès en  $O(1)$  en moyenne, quel que soit le nombre d'éléments dans la table. Toutefois, comme plusieurs données peuvent se trouver dans une même case, le temps d'accès dans le pire des cas peut être de  $O(n)$ . Comparées aux autres tableaux associatifs, les tables de hachage sont surtout utiles lorsque le nombre d'entrées est très important

7. Est-ce que le facteur de charge d'une table de Hashage. Est-ce que le contrôle du facteur de charge est nécessaire/optionnel pour le bon fonctionnement d'une table de Hashage avec Linear Probing ou Separate Chaining ? Quelle est la stratégie utilisée par `java.util.Hashtable` pour contrôler le facteur de charge ? En quoi est-elle différente de celle proposée dans `LinearProbinHashST` ? Quel est le lien entre le facteur de charge et collision ? (Sundeeep)

Le facteur de charge d'une table de hachage est le rapport entre le nombre d'éléments présents dans cette table et le nombre de places libres. Il y a un lien direct entre la notion de collision et celle du facteur de charge. En effet, plus le facteur de charge augmente, plus le tableau est rempli et plus le risque de collision augmente, ce qui affectera donc directement les performances du programme.

Cela signifie donc que, quelle que soit la manière de gérer les collisions (linear probing ou separate chaining), il est nécessaire de bien contrôler le facteur de charge.

Cependant, nous pourrions nous permettre un facteur de charge plus important dans le cas d'un separate chaining que dans le cas d'un linear probing, puisque le linear probing ne fonctionne plus du tout quand le load ratio atteint 1, et pour des valeurs proches de très longs clusters peuvent se former.

La stratégie utilisée dans `java.util.Hashtable` est de fixer le facteur de charge et la taille initiale de la table de hachage à sa création. À chaque fois que le facteur de charge est atteint, la nouvelle taille de la table est calculée en doublant la taille précédente.

La différence principale entre les stratégies de redimensionnement de `Hashtable` et `LinearProbingHashST` est que cette dernière redimensionne vers une taille plus petite quand la charge du tableau devient plus petite (quand le tableau est seulement un huitième plein). De cette manière, elle peut assurer que l'espace occupé par la table sera toujours proportionnel à son contenu.

8. Imaginez une nouvelle méthode `iterator()` qui retourne un itérateur sur les clefs de `LinearProbingHashST`. Votre itérateur ne devrait pas accepter de modification de la table de hashage alors qu'il est utilisée : une `ConcurrentModificationException()` doit être lancée si c'est le cas. Que suggérez-vous pour ce faire? Hint : Inspirez vous de la stratégie de `java.util.Hashtable`. (Thomas)

`Iterator()` parcourt la table de Hashage des clés jusqu'à M, en retournant à chaque fois la clés attendu grâce a la méthode `next()`. Il est bien évident qu'il contiendra les méthodes `hasnext()` qui lui indiquera si il y a encore un élément à lire après. Toutefois, si lors de l'appel à la méthode `next()`, le flag de modification est mise à vrai, `Iterator()` enverra `ConcurrentModificationException()`.

- 9.
- 10.
- 11.