

Soit E un ensemble d'éléments ayant chacun une clé. On note U l'ensemble des clés possibles. Quelles structures de données utiliser pour rechercher efficacement dans E ?

I Tableaux et listes

1) Vecteurs de bits ($U \subset \mathbb{N}$)

On représente E par un tableau T de taille $|U|$ tel que $T[i] = \begin{cases} 1 & \text{si } \exists x \in E, \text{clé}(x) = i \\ 0 & \text{sinon} \end{cases}$. L'avantage si $x \in E$ revient à savoir si $T[\text{clé}(x)] = 1$. On a donc une recherche en $O(1)$ mais cette idée est inutilisable si $|U|$ est grand.

2) Recherche naïve

On représente E par la liste L des clés de ses éléments. Rechercher x dans E revient à parcourir L jusqu'à trouver $\text{clé}(x)$. D'où une complexité au pire en $O(|E|)$.

3) Recherche dichotomique (U ordonné)

On représente E par le tableau trié de ses clés. On effectue alors une recherche dichotomique en $O(\log |E|)$.

II Hachage [Cor]

1) Principe général

Si $|U|$ est grand, I.1 est inutilisable. On cherche ici une structure de taille $O(|E|)$ telle que la recherche soit $O(1)$.

On choisit un tableau T de taille $m < |U|$ et $h: U \rightarrow \llbracket 0, m-1 \rrbracket$. On stocke alors $x \in E$ à la place $T[h(\text{clé}(x))]$. Comme $m < |U|$, h n'est pas injective. Si deux clés ont même image par h , on parle de collision (Figure 1).

2) Résolution par chaînage

On considère T comme un tableau de listes. Ainsi $T[i]$ contient tous les éléments $x \in E$ tel que $h(\text{clé}(x)) = i$. Notons alors $n_i = |T[i]|$ et $\alpha = |E|/m$ le coefficient de remplissage.

Si h est mal choisie, toutes les clés peuvent être envoyées sur la même case. Dans ce cas la recherche se fait en $O(|E|)$. On va voir comment choisir h pour éviter ça.

Hypothèse: On suppose que h vérifie l'hypothèse de hachage uniforme simple (HHUS):

$$P(h(x) = a) = 1/m \quad \forall x \in U \quad \forall a \in \llbracket 0, m-1 \rrbracket.$$

thm: Sous HHUS, toute recherche infructueuse est au pire en $O(1 + \alpha)$ et toute recherche réussie est en moyenne en $O(1 + \alpha)$.

3) Hachage universel et hachage parfait

def: Soit \mathcal{H} une classe de fonctions de hachage. \mathcal{H} est dite universelle si

$$\forall k \neq l \in U, |\{h \in \mathcal{H}, h(k) = h(l)\}| \leq |\mathcal{H}|/m.$$

ex: Si $p > |U|$ est premier, la classe $\{h: k \mapsto (ak + b) \bmod p \bmod m, a \in (\mathbb{Z}/p\mathbb{Z})^*, b \in \mathbb{Z}/p\mathbb{Z}\}$ est universelle.

thm: Si \mathcal{H} est universelle et $h \in \mathcal{H}$, $E(m_{h(R)}) \leq 1 + \alpha \quad \forall R \in U$.

app: Soit E un ensemble fini, i.e. on ne fait plus d'insertion ni DEV de suppression, seulement des recherches. On peut construire une structure de taille $O(|E|)$ où la recherche est $O(1)$ au pire.

app: Implémentation des dictionnaires

- utilisation de CD-RAM
- tables des symboles / mots réservés dans les compilateurs

III Structures arborescentes [Cor][Bear]

1) Arbres de recherche

def: Un arbre de recherche est un arbre étiqueté tel que tout nœud ayant $k+1$ sous-arbres $A_0 \dots A_k$ contient un k -uplet $x_1 < \dots < x_k$ tel que pour toute composante x d'une étiquette d'un nœud de A_i , $x_i < x < x_{i+1}$.

ex: Figure 2.

def: Un arbre linéaire de recherche (ABR) est un arbre de recherche où tout nœud non terminal a deux fils.

On a un algorithme de recherche linéaire en la hauteur de l'arbre, donc entre $\log |E|$ et $|E|$.

Si pour chaque $d \in T$, on connaît la probabilité qu'on la recherche, on peut s'intéresser à des ABR optimaux, i.e. qui minimisent l'espérance du temps de recherche d'une $d \in T$ (DEV2).

2) Équilibrage et AVL

def: Un AVL est un arbre binaire de recherche tel que pour tout nœud $x \in T$ de sous-arbres G et D , $|\ln(G) - \ln(D)| \leq 1$.

thm: Soit T un AVL à n nœuds, de hauteur h . Alors

$$\log_2(n+1) \leq h+1 \leq 1.44 \log_2(n+1)$$

cor: Les opérations de recherche dans un AVL sont en $O(\log(n))$.

remq: Les opérations d'insertion et suppression sont plus complexes à cause des rééquilibrages mais restent $O(\log n)$.

3) Arbres 2-3-4

def: Un arbre 2-3-4 est un arbre de recherche où tous les nœuds internes ont 2, 3 ou 4 fils, et dont toutes les feuilles sont à la même hauteur.

thm: La hauteur h d'un arbre 2-3-4 contenant n clés vérifie $\log_4(n+1) \leq h+1 \leq \log_2(n+1)$.

cor: La recherche est encore en $O(\log n)$.

4) Arbres rouge-noir

def: Un arbre rouge/noir est un ABR vérifiant:

- * tout nœud est rouge ou noir
- * la racine est noire
- * les feuilles sont noires
- * les fils d'un nœud rouge sont noirs
- * $\forall x$, les chemins issus de x et allant jusqu'aux feuilles ont même nombre de nœuds noirs.

thm: Un arbre rouge/noir à n nœuds a une hauteur plus petite que $2 \log_2(n+1)$.

cor: La recherche est en $O(\log n)$.

remq: On peut passer d'un arbre rouge/noir à un arbre 2-3-4 et réciproquement voir la figure 3.

5) Arbres persistants et duplication de chemins

On veut gérer une base de données qui évolue dans le temps. On veut savoir si x appartient à la structure au temps t (i.e. après t opérations insérer/supprimer).

- 1^{ère} idée: recopiage intégral à chaque opération (Figure 4)
2^{ème} idée: on utilise des arbres et on duplique chaque nœud dont un descendant est modifié (Figure 5).

La recherche de x à l'instant t revient à chercher x dans l'arbre dont la racine est pointée par $T[t]$. L'insertion, la suppression et la recherche sont en $O(\ln n)$ où n est le plus grand nombre de clés présentes dans la structure à l'instant t .

6) Tas et files de priorité

def: • Un tas est un arbre binaire dont tous les niveaux sont complets sauf éventuellement le dernier, où tous les nœuds sont tassés à gauche.

• Un tas max est un tas où chaque nœud est plus grand que ses fils.

smg: • On représente un tas par un tableau (Figure 6)

• Le maximum d'un tas max est sa racine.

Si A est un tas, $x \in A$ dont les sous-arbres G et D sont des tas max. Alors on peut obtenir un tas max contenant x, G et D en faisant descendre x tant qu'il viole la propriété de tas max. On parle d'entassement.

app: • Construction d'un tas max à partir d'un tas en temps linéaire

• On utilise l'entassement pour récupérer un tas max quand on extrait la racine d'un tas max.

On peut implémenter une file de priorité grâce aux tas max:

• défilage par extraction du maximum: on extrait la racine puis on entasse pour obtenir un tas max $\rightarrow O(\ln n)$

• enfilage: on ajoute un nœud (correspondant à une case à la fin du tableau) qui contient la nouvelle priorité et on fait remonter le nœud jusqu'à obtenir un tas max $\rightarrow O(\ln n)$.

app: • algorithme de Dijkstra

• ordonnancement de tâches/processus.

• codage de Huffman.

[Coi][Beau][Pop][Baa]

IV Union Find, union par rang et compression de chemins

On veut pouvoir gérer des classes d'équivalences de $\{0, n-1\}$ avec deux opérations: • déterminer la classe d'un élément
• fusionner deux classes.

On représente chaque classe par un arbre et la partition par une forêt. On prend un tableau T tel que $T[i]$ est le père de i . Si i est racine, on pose par convention $T[i] = i$.
On a alors les deux algorithmes

Trouver(x)

┌ Si $x \neq T[x]$, Trouver($T[x]$)
└ Sinon renvoyer x

Union(x, y)

┌ $T[\text{Trouver}(x)] \leftarrow \text{Trouver}(y)$

problème: On peut se retrouver avec des arbres peigne
Trouver est alors en $O(n)$

Union par rang:

→ On calcule un tableau Rang tel que Rang[x] majore la hauteur du sous-arbre issu de x .

→ Lors d'une union, on fait pointer le nœud de plus petit rang parmi Trouver(x) et Trouver(y) vers l'autre

→ On actualise Rang

Compression de chemins:

→ Quand on appelle Trouver(x), on fait pointer tous les nœuds parcourus vers Trouver(x)

(Figure 7)

flm: Si on fait m opérations, leur complexité est $O(m \alpha(m))$ donc la complexité amortie est $O(\alpha(m))$ où α est une fonction à croissance très lente ($\alpha(m) \leq 5$ pour tous les m utilisés en pratique).

app: • algorithme de Kruskal

* composantes connexes d'un graphe

* multiplication

* test de congruence.

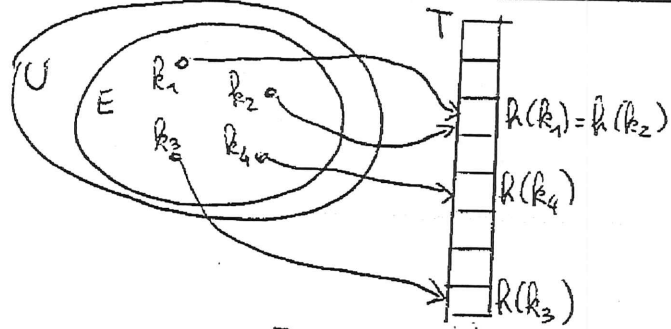


Figure 1.

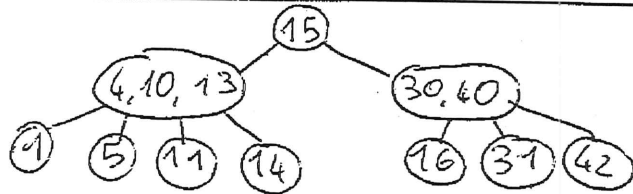


Figure 2

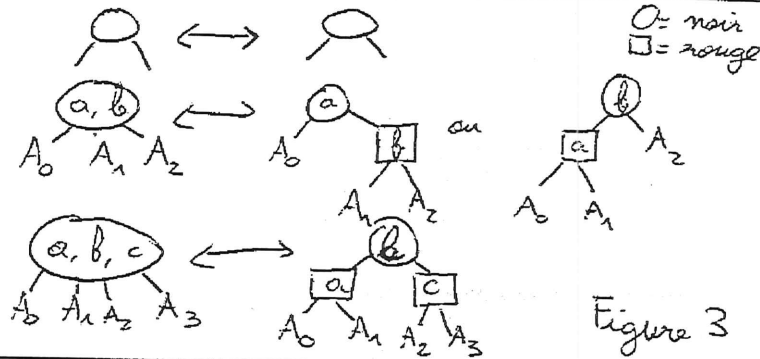


Figure 3

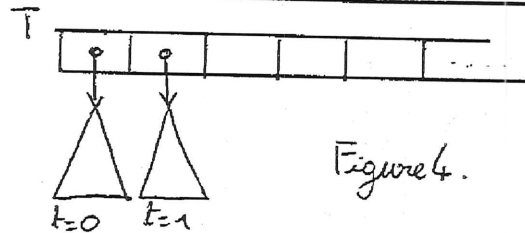


Figure 4.

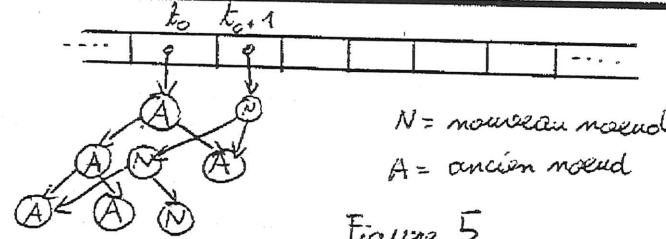


Figure 5

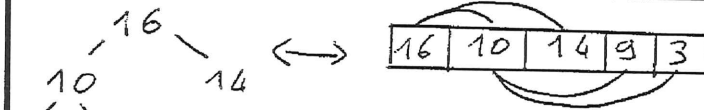


Figure 6

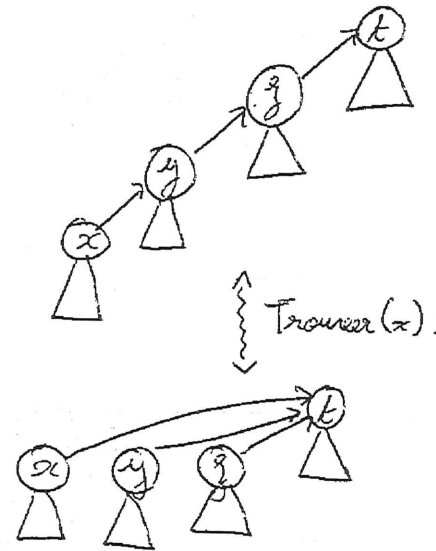


Figure 7.

[Cor] Cormen

[Beau] Beauquier

[Pap] Papadimitriou, Algorithms

[Boc] Bocca, Nephew, Term rewriting systems and all that