

- HAI503I: Algorithmique 4 -

Chap. 4 – Tables de hachage

L3 informatique
Université de Montpellier

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 Résolution par chaînage
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 Résolution par chaînage
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

1. Introduction

1.1 Structure de données dictionnaire

1.2 Tables de hachage

1.3 Fonctions de hachage

2. Résolution des collisions

2.1 Problématique

2.2 Résolution par chaînage

2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Les dictionnaires Python

Comment implanter le type dict de Python ?

```
1 >>> NbPattes = {}                                # Dict vide
2 >>> NbPattes['Humain'] = 2                        # Ajout
3 >>> NbPattes['Mille_Pattes'] = 999
4 >>> NbPattes['Araignee'] = 8
5 >>> NbPattes['Mille_Pattes'] = 1000            # Modification
6 >>> 'Coccinelle' in NbPattes                    # Recherche
7 False
8 >>> NbPattes['Mille_Pattes']
9 1000
```

La structure de données *dictionnaire*

Définition d'un dictionnaire

- ▶ Ensemble de couples (clef, valeur)
- ▶ Opérations disponibles :
 - ▶ **Création** d'un dictionnaire vide
 - ▶ **Insertion** d'un couple
 - ▶ Modification d'une valeur \rightsquigarrow Ré-Insertion
 - ▶ **Recherche** d'une clef \rightsquigarrow renvoie la valeur ou une erreur

Objectif

Les opérations RECHERCHE et INSERTION doivent être *rapides*

La structure de données *dictionnaire*

Définition d'un dictionnaire

- ▶ Ensemble de couples (clef, valeur)
- ▶ Opérations disponibles :
 - ▶ **Création** d'un dictionnaire vide
 - ▶ **Insertion** d'un couple
 - ▶ Modification d'une valeur \rightsquigarrow Ré-Insertion
 - ▶ **Recherche** d'une clef \rightsquigarrow renvoie la valeur ou une erreur

Objectif

Les opérations RECHERCHE et INSERTION doivent être *rapides*

Hypothèse simplificatrice

- ▶ Les clefs sont des entiers
- ▶ Théorie : toute donnée est codée en binaire \rightsquigarrow interprétation comme un entier
- ▶ Pratique : codage ASCII, Unicode ou transformation quelconque...

Quelles solutions ?

Dictionnaire de n éléments, clef entre 0 et $N - 1$ *

Tableau

- ▶ Taille : N ✗
- ▶ CRÉATION : $O(N)$ ✗
- ▶ INSERTION : $O(1)$ ✓
- ▶ RECHERCHE : $O(1)$ ✓

. * Dans l'exemple de départ : mots de 20 lettres sur un alphabet de taille 128 \rightsquigarrow
 $N = 128^{20} = 2^{140} \sim 1.4 \times 10^{42} \dots$

Quelles solutions ?

Dictionnaire de n éléments, clef entre 0 et $N - 1^*$

Tableau

- ▶ Taille : N ✗
- ▶ CRÉATION : $O(N)$ ✗
- ▶ INSERTION : $O(1)$ ✓
- ▶ RECHERCHE : $O(1)$ ✓

Liste chaînée

- ▶ Taille : n ✓
- ▶ CRÉATION : $O(1)$ ✓
- ▶ INSERTION : $O(n)$ ✗
- ▶ RECHERCHE : $O(n)$ ✗

. * Dans l'exemple de départ : mots de 20 lettres sur un alphabet de taille 128 \rightsquigarrow
 $N = 128^{20} = 2^{140} \sim 1.4 \times 10^{42} \dots$

Quelles solutions ?

Dictionnaire de n éléments, clef entre 0 et $N - 1^*$

Tableau

- ▶ Taille : N ❌
- ▶ CRÉATION : $O(N)$ ❌
- ▶ INSERTION : $O(1)$ ✓
- ▶ RECHERCHE : $O(1)$ ✓

Arbre binaire de recherche

- ▶ Taille : n ✓
- ▶ CRÉATION : $O(1)$ ✓
- ▶ INSERTION : $O(h) \rightsquigarrow O(\log n)$ si équilibré 😊
- ▶ RECHERCHE : $O(h) \rightsquigarrow O(\log n)$ si équilibré 😊

Liste chaînée

- ▶ Taille : n ✓
- ▶ CRÉATION : $O(1)$ ✓
- ▶ INSERTION : $O(n)$ ❌
- ▶ RECHERCHE : $O(n)$ ❌

. * Dans l'exemple de départ : mots de 20 lettres sur un alphabet de taille 128 \rightsquigarrow
 $N = 128^{20} = 2^{140} \sim 1.4 \times 10^{42} \dots$

Quelles solutions ?

Dictionnaire de n éléments, clef entre 0 et $N - 1^*$

Tableau

- ▶ Taille : N ❌
- ▶ CRÉATION : $O(N)$ ❌
- ▶ INSERTION : $O(1)$ ✅
- ▶ RECHERCHE : $O(1)$ ✅

Arbre binaire de recherche

- ▶ Taille : n ✅
- ▶ CRÉATION : $O(1)$ ✅
- ▶ INSERTION : $O(h) \rightsquigarrow O(\log n)$ si équilibré 😊
- ▶ RECHERCHE : $O(h) \rightsquigarrow O(\log n)$ si équilibré 😊

Liste chaînée

- ▶ Taille : n ✅
- ▶ CRÉATION : $O(1)$ ✅
- ▶ INSERTION : $O(n)$ ❌
- ▶ RECHERCHE : $O(n)$ ❌

Tas

- ▶ Taille : n ✅
- ▶ CRÉATION : $O(1)$ ✅
- ▶ INSERTION : $O(\log n)$ 😊
- ▶ RECHERCHE : $O(\log n)$ 😊

. * Dans l'exemple de départ : mots de 20 lettres sur un alphabet de taille 128 \rightsquigarrow
 $N = 128^{20} = 2^{140} \sim 1.4 \times 10^{42} \dots$

1. Introduction

1.1 Structure de données dictionnaire

1.2 Tables de hachage

1.3 Fonctions de hachage

2. Résolution des collisions

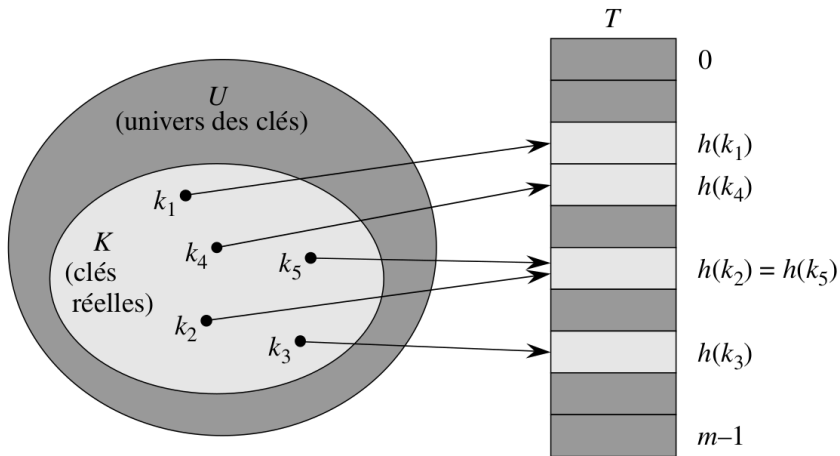
2.1 Problématique

2.2 Résolution par chaînage

2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Tables de hachage



Formalisation

Table de hachage

▶ Clefs

- ▶ Univers U des clefs possibles : $U = \{0, \dots, N - 1\}$
- ▶ Clefs utilisées : $K \subset U$, de taille n

▶ Table T de taille m

- ▶ Indices entre 0 et $m - 1$
- ▶ Une case peut être vide ou contenir une valeur, voire plusieurs

▶ Fonction de hachage

- ▶ Fonction $h : U \rightarrow \{0, \dots, m - 1\}$

Insertion

Insertion du couple (k, v) dans la case $T_{[h(k)]}$

Utilisées de partout :

- ▶ En machine : cache, compilateur (variables), disque (table des inodes)
- ▶ Pour gérer des bases de données
- ▶ En cryptographie, pour la vérification d'intégrité d'un fichier (par exemple, *la fonction MD5*), pour la vérification de mots de passe...

Questions à résoudre

Caractéristiques

- ▶ Taille : $m \rightsquigarrow$ comment choisir m par rapport à n et N ?
- ▶ CRÉATION : on veut $O(m)$
On ne peut pas explicitement fixer une valeur $h(k)$ pour chaque clef k de U (il y en a trop...), il va falloir une formule ou un algorithme pour calculer $h(k)$!
- ▶ INSERTION : calcul de $h(k)$ puis insertion en case $h(k) \rightsquigarrow$ quelle complexité?
- ▶ RECHERCHE : calcul de $h(k)$ puis recherche dans la case $h(k) \rightsquigarrow$ quelle complexité?

Questions à résoudre

Caractéristiques

- ▶ Taille : $m \rightsquigarrow$ comment choisir m par rapport à n et N ?
- ▶ CRÉATION : on veut $O(m)$
On ne peut pas explicitement fixer une valeur $h(k)$ pour chaque clef k de U (il y en a trop...), il va falloir une formule ou un algorithme pour calculer $h(k)$!
- ▶ INSERTION : calcul de $h(k)$ puis insertion en case $h(k) \rightsquigarrow$ quelle complexité ?
- ▶ RECHERCHE : calcul de $h(k)$ puis recherche dans la case $h(k) \rightsquigarrow$ quelle complexité ?

Collisions

- ▶ Que fait-on si $h(k_1) = h(k_2)$?
 - ▶ Plusieurs valeurs dans une case (liste chaînée, etc.)
 - ▶ Utiliser une autre case ?
- ▶ Est-ce que $h(k_1) = h(k_2)$ arrive souvent ?
 - ▶ Comment choisir h ?

1. Introduction

1.1 Structure de données dictionnaire

1.2 Tables de hachage

1.3 Fonctions de hachage

2. Résolution des collisions

2.1 Problématique

2.2 Résolution par chaînage

2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Problématique des fonctions de hachage

Contexte

- ▶ Choix d'une fonction $h : U = \{0, \dots, N - 1\} \rightarrow \{0, \dots, m - 1\}$
- ▶ Fonction utilisée pour un ensemble de clefs K de taille $n \ll N$

Quelques exemples (qui ne marchent pas forcément bien...)

- ▷ $h(i) = 0$ (pour tout $i \in U$)
- ▷ $h(i) = i \bmod m$
- ▷ $h(i) = a.i + b \bmod m$ (pour a et b à choisir)
- ▷ $h(i) = 63.i^2 - 12.i + 7 \bmod m$ (par exemple...)

Collisions évitables ?

- ▶ Avec $N \gg m$, on aura forcément des collisions ($h(k_1) = h(k_2)$)
- ▶ Mais on stocke n clefs avec $n \leq m$.
 - ▶ Pour un ensemble de clefs, possible de trouver h sans collision
 - ▶ Mais... on ne connaît pas les clefs à l'avance

Problématique des fonctions de hachage

Contexte

- ▶ Choix d'une fonction $h : U = \{0, \dots, N - 1\} \rightarrow \{0, \dots, m - 1\}$
- ▶ Fonction utilisée pour un ensemble de clefs K de taille $n \ll N$

Quelques exemples (qui ne marchent pas forcément bien...)

- ▷ $h(i) = 0$ (pour tout $i \in U$)
- ▷ $h(i) = i \bmod m$
- ▷ $h(i) = a.i + b \bmod m$ (pour a et b à choisir)
- ▷ $h(i) = 63.i^2 - 12.i + 7 \bmod m$ (par exemple...)

Problématique

- ▶ On veut choisir h avant de connaître les clefs
- ▶ On voudrait éviter les collisions entre clefs... sans les connaître
- ▶ **Pas le choix : une fonction de hachage doit être choisie aléatoirement**

Modèles aléatoires des fonctions de hachage

Tentative de choix de fonctions de hachage

On tire h uniformément parmi les fonctions de $U = \{0, \dots, N - 1\}$ dans $\{0, \dots, m - 1\}$

Représentation de h

- ▶ Pour chaque $k \in U$, une valeur $h(k) \rightsquigarrow$ tableau H de taille N
- ▶ Tirage de $h \rightsquigarrow$ tirage uniforme et indépendant de chaque $h(k)$ dans $\{0, \dots, m - 1\}$

Avantage et inconvénient

- ▶ *Avantage* : on a bien pour tout $k_1 \neq k_2$, $\Pr[h(k_1) = h(k_2)] = 1/m$
- ▶ *Inconvénient* : **totallement irréaliste** \rightsquigarrow **tableau de taille $N!!!$**

Modèles aléatoires des fonctions de hachage

Tentative de choix de fonctions de hachage

On tire h uniformément parmi les fonctions de $U = \{0, \dots, N - 1\}$ dans $\{0, \dots, m - 1\}$

Représentation de h

- ▶ Pour chaque $k \in U$, une valeur $h(k) \rightsquigarrow$ tableau H de taille N
- ▶ Tirage de $h \rightsquigarrow$ tirage uniforme et indépendant de chaque $h(k)$ dans $\{0, \dots, m - 1\}$

Avantage et inconvénient

- ▶ *Avantage* : on a bien pour tout $k_1 \neq k_2$, $\Pr[h(k_1) = h(k_2)] = 1/m$
- ▶ *Inconvénient* : **totallement irréaliste** \rightsquigarrow **tableau de taille $N!!!$**

Nécessité d'un autre modèle

- ▶ Avec les mêmes (bonnes) propriétés probabilistes
- ▶ Mais, avec qui soit implémentable en temps raisonnable

Modèle universel des fonctions de hachage

Tentative 2 de choix de fonctions de hachage

On fixe un ensemble \mathcal{H} de fonctions de hachage et on tire h uniformément dans \mathcal{H} .

Définition

Un ensemble \mathcal{H} de fonctions de $U = \{0, \dots, N - 1\}$ dans $\{0, \dots, m - 1\}$ est **universel** si pour tout $k_1 \neq k_2 \in U$, lorsqu'on tire uniformément h dans \mathcal{H} on a $\Pr[h(k_1) = h(k_2)] \leq 1/m$.

Remarques

- ▶ Probabilité de collision \leq probabilité dans le modèle aléatoire
- ▶ On *connait* des ensembles \mathcal{H} universels réalistes (cf Partie 3...)

Ensemble universel *intéressant*

- ▶ Ensemble \mathcal{H} pas trop gros \rightsquigarrow représentation de h assez petite
- ▶ Tirer uniformément $h \in \mathcal{H}$ doit être efficace
- ▶ Calculer $h(k)$ doit être rapide

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 Résolution par chaînage
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 **Problématique**
- 2.2 Résolution par chaînage
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Collisions

Contexte

- ▶ Table T avec fonction de hachage
 $h : U = \{0, \dots, N - 1\} \rightarrow \{0, \dots, m - 1\}$ choisie dans un ensemble universel
- ▶ Ensemble de clefs $K \subseteq U$
- ▶ Il y a **collision** entre les clefs k_1 et k_2 si $h(k_1) = h(k_2)$

Remarque

- ▶ Si la table est suffisamment grande, avec bonne probabilité, il n'y aura pas de collision.
- ▶ Mais on ne peut pas être complètement sûr de les éviter...

Collisions

Contexte

- ▶ Table T avec fonction de hachage $h : U = \{0, \dots, N - 1\} \rightarrow \{0, \dots, m - 1\}$ choisie dans un ensemble universel
- ▶ Ensemble de clefs $K \subseteq U$
- ▶ Il y a **collision** entre les clefs k_1 et k_2 si $h(k_1) = h(k_2)$

Remarque

- ▶ Si la table est suffisamment grande, avec bonne probabilité, il n'y aura pas de collision.
- ▶ Mais on ne peut pas être complètement sûr de les éviter...

Lemme

Si $m = n^2$, et h est tirée uniformément dans un ensemble universel \mathcal{H} , alors la probabilité qu'il existe deux clefs $k_1 \neq k_2$ telles que $h(k_1) = h(k_2)$ est $\leq \frac{1}{2}$.

Collisions

Contexte

- ▶ Table T avec fonction de hachage $h : U = \{0, \dots, N - 1\} \rightarrow \{0, \dots, m - 1\}$ choisie dans un ensemble universel
- ▶ Ensemble de clefs $K \subseteq U$
- ▶ Il y a **collision** entre les clefs k_1 et k_2 si $h(k_1) = h(k_2)$

Remarque

- ▶ Si la table est suffisamment grande, avec bonne probabilité, il n'y aura pas de collision.
- ▶ Mais on ne peut pas être complètement sûr de les éviter...

Deux exemples de (familles de) solutions

- ▶ Plusieurs éléments dans une même case : *résolution par chaînage*
- ▶ Trouver une autre case libre : *adressage ouvert*

1. Introduction

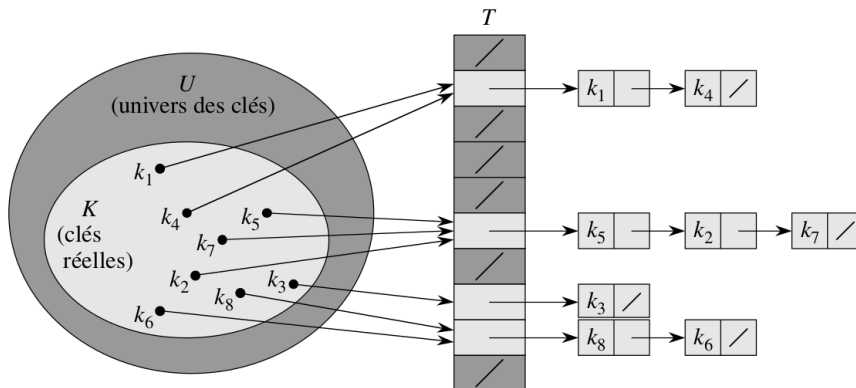
- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 **Résolution par chaînage**
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Résolution par chaînage : principe



Résolution par chaînage

Chaque case de T contient une liste chaînée

Algorithmes

- ▶ RECHERCHE de k :
 - ▶ Calcul de $h(k)$
 - ▶ Parcours de la liste contenue dans $T_{[h(k)]}$, retour de (k, v) ou NULL si k n'apparaît pas dans la liste.
 - ▶ Complexité : $O(\ell(k))$ où $\ell(k)$ est la taille de la liste $T_{[h(k)]}$
- ▶ INSERTION de (k, v) :
 - ▶ Si k apparaît déjà dans la liste contenue dans $T_{[h(k)]}$, on remplace sa valeur par v
 - ▶ Sinon, on ajoute (k, v) à la liste $T_{[h(k)]}$
 - ▶ Complexité : $O(\ell(k))$ où $\ell(k)$ est la taille de la liste $T_{[h(k)]}$

Quelle efficacité ?

- ▶ Une opération coûte $O(L)$, où $L = \max_{k \in K} \ell(k) \rightsquigarrow$ quelle taille maximale en moyenne ?

Efficacité de la résolution par chaînage

Théorème

Soit T une table de hachage de taille m , avec h tirée uniformément dans un ensemble \mathcal{H} universel. Si T contient n éléments et que les collisions sont résolues par chaînage, l'espérance de la complexité de RECHERCHE et de INSERTION est en $O(n/m)$.

Complexité

- ▶ Complexité espérée : $O(\alpha)$ où $\alpha = \frac{n}{m}$ est le *taux de remplissage*
- ▶ Si le taux est autour de 1 : $O(1)$ en moyenne
- ▶ La résolution par chaînage marche bien en moyenne, mais certaines opérations peuvent être coûteuses (si on a une grande liste...)

Pourquoi des listes chaînées ?

- ▶ Arbres binaires de recherche ou tas dans chaque case
 - ▶ Complexité moyenne en $O(\log \alpha)$
 - ▶ Complexité pire cas en $\max_k \log \ell(k)$
- ▶ Et pourquoi pas des tables de hachage ?

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 Résolution par chaînage
- 2.3 **Adressage ouvert**

3. Une famille universelle de fonctions de hachage

Principe

Idée

Si la case pour insérer (k, v) est occupée, trouver une autre case

Formellement

- ▶ m fonctions de hachage h_1, \dots, h_m
 - ▶ 1^{er} essai : INSERTION en case $h_1(k)$
 - ▶ 2^{ème} essai : INSERTION en case $h_2(k)$
 - ▶ ...
 - ▶ $m^{\text{ème}}$ essai : INSERTION en case $h_m(k)$
- ▶ Condition : pour tout k , $\{h_1(k), \dots, h_m(k)\}$ est une **permutation** de $\{0, \dots, m-1\}$

Algorithmes

- ▶ RECHERCHE : explorer $T_{[h_1(k)]}, T_{[h_2(k)]}, \dots$
 - ▶ si on trouve $k \rightsquigarrow$ gagné
 - ▶ si on trouve une case vide \rightsquigarrow k n'est pas dans T
- ▶ INSERTION : explorer jusqu'à trouver une case vide

Adressage ouvert : constructions et performance

Construire les m fonctions à partir d'une (ou deux) fonctions de hachage.

Quelques possibilités pratiques

- ▶ **Sondage linéaire** : $h_i(k) = (h(k) + i) \bmod m$
- ▶ **Sondage quadratique** : $h_i(k) = (h(k) + ai^2 + bi) \bmod m$ (a et b à choisir)
- ▶ **Sondage binaire** : $h_i(k) = h(k) \oplus i$ (si $m = 2^\ell$)
- ▶ **Double hachage** : $h_i(k) = (h^{(1)}(k) + ih^{(2)}(k)) \bmod m$ ($h^{(1)}$ et $h^{(2)}$ à choisir)
- ▶ ...

Hypothèse (théorique...)

Pour tout k , $\{h_1(k), \dots, h_m(k)\}$ est une permutation aléatoire

Théorème

Sous l'hypothèse précédente, si le facteur de remplissage est $\alpha = n/m < 1$, l'espérance du nombre de cases visitées pour une RECHERCHE infructueuse ou une INSERTION est $\leq \frac{1}{1-\alpha}$.

Bilan sur l'adressage ouvert

Principe de base

- ▶ Une seule table principale, un seul élément par case
- ▶ Si une case est occupée, *aller ailleurs*, plusieurs possibilités pour ça...

Complexité espérée (modèle aléatoire)	$\alpha = \frac{1}{2}$	$\alpha = \frac{9}{10}$
▶ INSERTION ou RECHERCHE infructueuse : $\frac{1}{1-\alpha}$	2	10
▶ RECHERCHE réussie : $\frac{1}{\alpha} \ln \frac{1}{1-\alpha}$ (admis)	$\leq 1,387$	$\leq 2,559$

Bilan sur l'adressage ouvert

Principe de base

- ▶ Une seule table principale, un seul élément par case
- ▶ Si une case est occupée, *aller ailleurs*, plusieurs possibilités pour ça...

Complexité espérée (modèle aléatoire)	$\alpha = \frac{1}{2}$	$\alpha = \frac{9}{10}$
▶ INSERTION ou RECHERCHE infructueuse : $\frac{1}{1-\alpha}$	2	10
▶ RECHERCHE réussie : $\frac{1}{\alpha} \ln \frac{1}{1-\alpha}$ (admis)	$\leq 1,387$	$\leq 2,559$

Un exemple supplémentaire : *hachage du coucou*

- ▶ Deux fonctions de hachage $h^{(1)}$ et $h^{(2)}$ *deux emplacements possibles par clef*
- ▶ INSERTION de (k, v) :
 - ▶ Insertion en case $h^{(1)}(k)$
 - ▶ Si la case contenait (k', v') , on déplace (k', v') dans la case $h^{(2)}(k')$
 - ▶ Et récursivement...
- ▶ Et ça marche !

Conclusion sur la résolution des collisions

Les collisions sont inévitables

Deux familles de résolutions vues ici

▶ **Chaînage :**

- ▶ Gérer les collisions en mettant plusieurs éléments par case
- ▶ Complexité liée au nombre maximal d'éléments par case et à la structure de données

▶ **Adressage ouvert :**

- ▶ Gérer les collisions en cherchant une autre case libre
- ▶ Complexité liée au nombre de cases à inspecter

↪ Dans les deux cas : **complexité liée au nombre de collisions**

Conclusion sur la résolution des collisions

Les collisions sont inévitables

Deux familles de résolutions vues ici

▶ **Chaînage :**

- ▶ Gérer les collisions en mettant plusieurs éléments par case
- ▶ Complexité liée au nombre maximal d'éléments par case et à la structure de données

▶ **Adressage ouvert :**

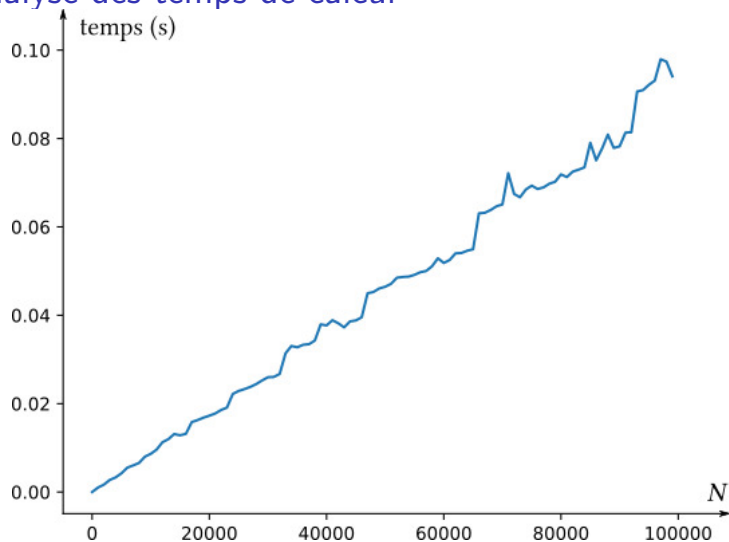
- ▶ Gérer les collisions en cherchant une autre case libre
- ▶ Complexité liée au nombre de cases à inspecter

↪ Dans les deux cas : **complexité liée au nombre de collisions**

Cas des dictionnaires Python

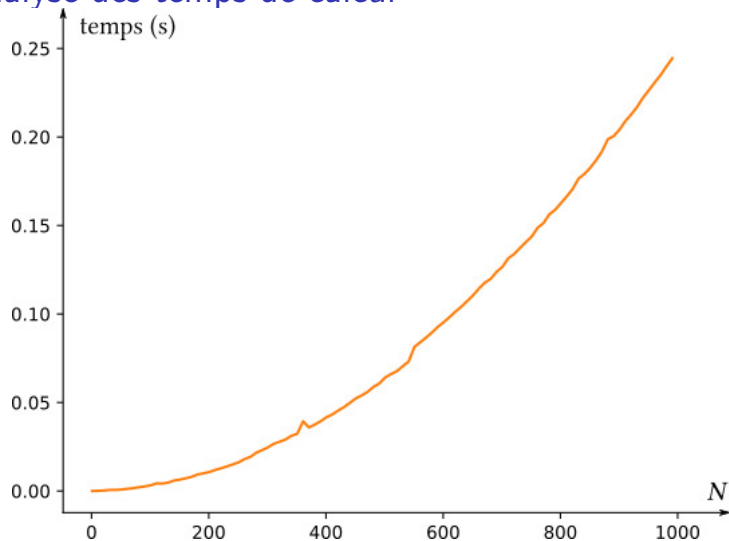
- ▶ Fonction de hachage pas aléatoire ↪ $h(i) = i \bmod (2^{61} - 1)$
- ▶ Résolution des collisions par adressage ouvert
 - ▶ Ordre de parcours des cases un peu complexe
- ▶ Solution théoriquement faible, à peu près correcte en pratique

Analyse des temps de calcul



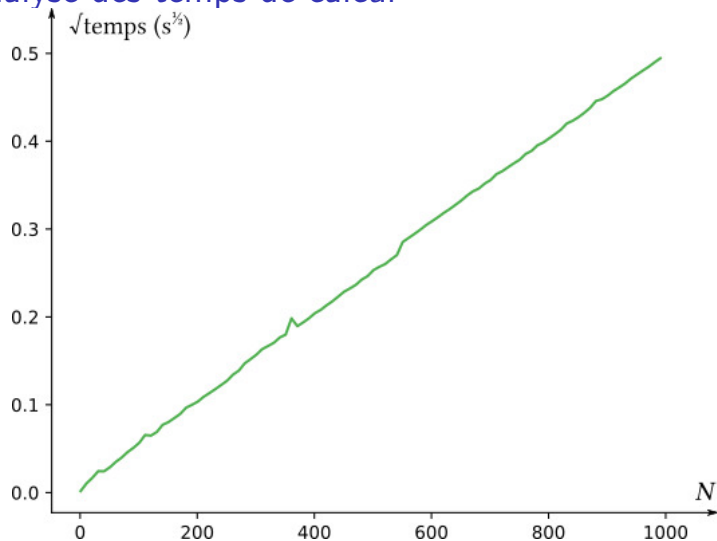
```
1 d={}
2 for i in range(N):
3     d[randrange(2**61*N**2)]=i
```

Analyse des temps de calcul



```
1 d={}
2 for i in range(N):
3     d[(2**61-1)*randrange(N**2)]=i
```


Analyse des temps de calcul



```
1 d={}
2 for i in range(N):
3     d[(2**61-1)*randrange(N**2)]=i
```

1. Introduction

- 1.1 Structure de données dictionnaire
- 1.2 Tables de hachage
- 1.3 Fonctions de hachage

2. Résolution des collisions

- 2.1 Problématique
- 2.2 Résolution par chaînage
- 2.3 Adressage ouvert

3. Une famille universelle de fonctions de hachage

Objectif

Rappel de la définition

Un ensemble \mathcal{H} de fonctions de $\{0, \dots, N-1\}$ dans $\{0, \dots, m-1\}$ est **universel** si pour tout $k_1 \neq k_2$ fixés, $\Pr[h(k_1) = h(k_2)] \leq 1/m$.

Contraintes sur \mathcal{H}

- ▶ Suffisamment grand pour avoir une probabilité $\leq 1/m$
- ▶ Suffisamment petit pour savoir représenter $h \in \mathcal{H}$ avec une place raisonnable
- ▶ Suffisamment *simple* pour savoir tirer $h \in \mathcal{H}$ en temps raisonnable

On veut \mathcal{H} de taille polynomiale en N

- ▶ Nombre de couples de clefs possibles $\binom{N}{2} \rightsquigarrow$ on demande au moins autant de fonctions h , donc **on veut** $|\mathcal{H}| \geq \binom{N}{2}$
- ▶ **Représentation d'une fonction h en $O(\log N)$ bits** \rightsquigarrow similaire à une clef
- ▶ **Tirage aléatoire en $O(\log N)$** \rightsquigarrow équivalent au calcul de $h(k)$

Hachage multiplicatif

Définition

Soit $\mathcal{H}_p^{N,m} = \{h_{a,b} : 0 < a < p, 0 \leq b < p\}$ la famille de fonctions définies par

$$h_{a,b} : \begin{cases} \{0, \dots, N-1\} & \rightarrow \{0, \dots, m-1\} \\ k & \rightarrow ((ak + b) \bmod p) \bmod m \end{cases}$$

où p est un nombre premier $> N$

Représentation et tirage aléatoire

- ▶ Tirage aléatoire de $h_{a,b}$: tirage de $a \in \{1, \dots, p-1\}$ et $b \in \{0, \dots, p-1\}$
- ▶ Représentation de $h_{a,b}$: (a, b, p)
- ▶ Taille : $|\mathcal{H}_p^{N,m}| = p(p-1) > N^2$

Hachage multiplicatif

Définition

Soit $\mathcal{H}_p^{N,m} = \{h_{a,b} : 0 < a < p, 0 \leq b < p\}$ la famille de fonctions définies par

$$h_{a,b} : \begin{cases} \{0, \dots, N-1\} & \rightarrow \{0, \dots, m-1\} \\ k & \rightarrow ((ak + b) \bmod p) \bmod m \end{cases}$$

où p est un nombre premier $> N$

Représentation et tirage aléatoire

- ▶ Tirage aléatoire de $h_{a,b}$: tirage de $a \in \{1, \dots, p-1\}$ et $b \in \{0, \dots, p-1\}$
- ▶ Représentation de $h_{a,b}$: (a, b, p)
- ▶ Taille : $|\mathcal{H}_p^{N,m}| = p(p-1) > N^2$

Théorème

La famille $\mathcal{H}_p^{N,m}$ est universelle (pour tout N, m et $p \geq N$)

Preuve du théorème

Rappel de la définition de $h_{a,b}$

Soit $\mathcal{H}_p^{N,m} = \{h_{a,b} : 0 < a < p, 0 \leq b < p\}$ la famille de fonctions définies par

$$h_{a,b} : \begin{array}{l} \{0, \dots, N-1\} \rightarrow \{0, \dots, m-1\} \\ k \rightarrow ((ak + b) \bmod p) \bmod m \end{array}$$

où p est un nombre premier $> N$

Lemme : système linéaire *modulo* p

Soit $k_1 \neq k_2$ et $u \neq v$ dans $\mathbb{Z}/p\mathbb{Z}$, alors il existe un unique couple $a, b \in \mathbb{Z}/p\mathbb{Z}$ tel que $u = ak_1 + b$ et $v = ak_2 + b$.

Preuve du théorème

Rappel de la définition de $h_{a,b}$

Soit $\mathcal{H}_p^{N,m} = \{h_{a,b} : 0 < a < p, 0 \leq b < p\}$ la famille de fonctions définies par

$$h_{a,b} : \begin{array}{l} \{0, \dots, N-1\} \rightarrow \{0, \dots, m-1\} \\ k \rightarrow ((ak + b) \bmod p) \bmod m \end{array}$$

où p est un nombre premier $> N$

Lemme : système linéaire *modulo* p

Soit $k_1 \neq k_2$ et $u \neq v$ dans $\mathbb{Z}/p\mathbb{Z}$, alors il existe un unique couple $a, b \in \mathbb{Z}/p\mathbb{Z}$ tel que $u = ak_1 + b$ et $v = ak_2 + b$.

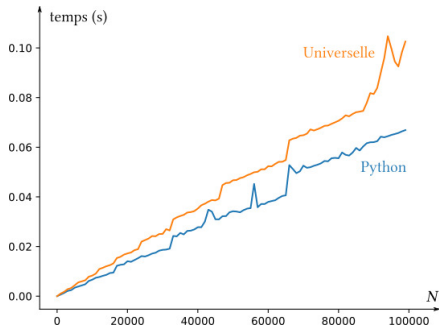
Théorème (réécrit)

Pour tout $k_1 \neq k_2$, $\Pr[h_{a,b}(k_1) = h_{a,b}(k_2)] \leq 1/m$

Bilan sur la famille universelle

Utilisation de la famille

- ▶ CRÉATION du dictionnaire : tirage aléatoire de a et b
- ▶ Complexité du calcul de $h_{a,b}(k) = ((ak + b) \bmod p) \bmod m$
 - ▶ Additions, multiplications, divisions d'entiers $\leq p^2$:
 $O(\log^2 p) = O(\log^2 N)$
 - ▶ Taille d'une clef $\rightsquigarrow O(\log N)$



Bilan sur la famille universelle

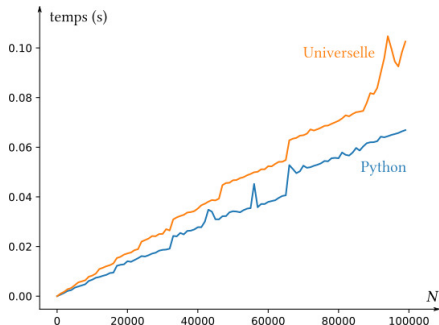
Utilisation de la famille

- ▶ CRÉATION du dictionnaire : tirage aléatoire de a et b
- ▶ Complexité du calcul de $h_{a,b}(k) = ((ak + b) \bmod p) \bmod m$
 - ▶ Additions, multiplications, divisions d'entiers $\leq p^2$:
 $O(\log^2 p) = O(\log^2 N)$
 - ▶ Taille d'une clef $\rightsquigarrow O(\log N)$

Autres familles universelles

- ▶ $h_a(k) = (ak \bmod 2^w) \text{ div } 2^{w-\ell}$
- ▶ $h_{\vec{c}}(k) = ((\sum_i c_i k^i) \bmod p) \bmod m$

quasi-universelle (cf TD)
fortement universelle



Conclusion sur les tables de hachage

Tables de hachage

- ▶ Structure de données très efficace, et très répandue
- ▶ Autres structures dérivées (ex : filtres de Bloom)
- ▶ Constructions pratiques inspirées de la théorie

Gestion des collisions

- ▶ Chaînage \rightsquigarrow complexité amortie $O(1)$ dans le modèle universel
- ▶ Adressage ouvert \rightsquigarrow complexité amortie $O(1)$ dans le modèle aléatoire
- ▶ D'autres méthodes existent. Ex. le hachage parfait \rightsquigarrow complexité pire cas $O(1)$ dans le modèle universel

Construction de familles universelles

- ▶ $h_{a,b}(k) = (((ak + b) \bmod p) \bmod m)$ fournit une famille universelle
- ▶ D'autres familles existent...

Dans les langages de programmation

- ▶ Tables de hachages souvent proposées (dictionnaires), non aléatoires
- ▶ Souvent bon en pratique, mais mauvaises surprises possibles