

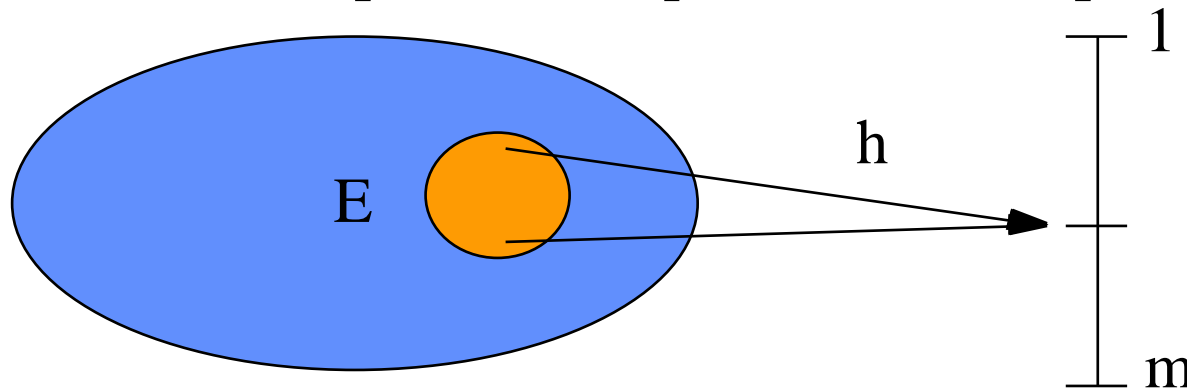
# Le hachage

# Principe général (1)

- Arbres de recherche => élément placé par rapport aux autres
- Hachage => élément placé à partir d'un calcul sur sa clé
  - complexité moyenne constante (indépendant du nombre d'éléments)
  - adaptation facile sur disque
- Remarque: cardinalité  $\{c \mid c \in E\}$  très inférieur au nombre de clés différentes:
  - 500 000 personnes, clés sur 30 octets =>  $3 \cdot 10^{42}$  clés
  - algorithme de placement non injectif

# Principe général (2)

- $h$  : clé  $[1..m]$
- tableau  $T [1..m]$
- si pour tout  $c, c' \in E, h(c) = h(c')$ , placer l'élément de clé  $c$  en  $h(c) \Rightarrow T[h(\text{la\_clé}(e))] = e$
- sinon  $\Rightarrow$  collisions,
- $h$  localise une partie de  $E$  qui est haché en petits morceaux



# Principe général (3)

- Il n'est pas question d'éliminer les collisions  
groupe de 23 personnes, la probabilité que deux personnes aient leur anniversaire le même jour est de 0,5  
en prenant  $m = 365$ , placer 23 personnes  $\Rightarrow$  1 chance sur 2 collision
- choix judicieux de  $h$  limite les collisions  
probabilité pour que  $\#\{c \mid c \in E \text{ et } h(c) = i\} = k$

	0	1	2	3	4	5
#E = 0,5m	0,60	0,30	0,08	0,01	0,002	0,0002
#E = m	0,37	0,37	0,18	0,06	0,01	0,003

# Spécification du hachage

**type** ens-hach **dérive de** table [ensemble [élément, clé]]

**paramètre** élément, clé **utilise** booléen

**avec** la\_clé: élément clé

h: clé 1 .. m

**opérations**

vide-h : ens-hach

ajout-h : élément × ens-hach / ens-hach

supprim-h : clé × ens-hach / ens-hach

recherch-h : clé × ens-hach / élément

$\_ \_ h \_$  : clé × ens-hach booléen

**préconditions**

$\neg(\text{la\_clé}(e) \_ h E)$

$c \_ h E$

$c \_ h E$

**axiomes**

$\text{borneinf}(\text{vide-h}) = 1 \quad \text{bornesup}(\text{vide-h}) = m$  ]

$1 \leq i \leq m \quad \text{ième}(\text{vide-h}, i) = \text{vide}$

$c \_ h E \quad c \quad \text{ième}(E, h(c))$

$\text{ajout-h}(e, E) = \text{soit } i = h(\text{la\_clé}(e)); \text{changer-ième}(E, i, \text{ajouter}(e, \text{ième}(E, i)))$

$\text{supprim-h}(c, E) = \text{soit } i = h(c); \text{changer-ième}(E, i, \text{supprimer}(c, \text{ième}(E, i)))$

$\text{recherch-h}(c, E) = \text{recherche}(c, \text{ième}(E, h(c)))$

*initialement table de m ensembles vides*

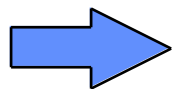


# Fonction de hachage

- Elle doit être:
  - équiprobable  $\Rightarrow$  probabilité que  $h(c) = i$  est  $1/m$
  - déterministe  $\Rightarrow$  pour  $c$  donnée,  $h(c)$  donne toujours la même valeur
  - facilement calculable  $\Rightarrow$  bonnes performances
- Idées générales
  - extraction de bits  $\Rightarrow$  extraire  $p$  bits de la clé, avec  $m = 2^p$
  - compression de bits  $\Rightarrow$  découpage de la clé en  $q$  tranches de  $p$  bits puis combiner les tranches par xor et rotations
  - division  $\Rightarrow h(c) = c \bmod m$
  - multiplication  $\Rightarrow h(c) = ((c * r) \bmod 1) * m$ .
- Pas de fonction universelle  $\Rightarrow$  essayer sur échantillon

# Résolution des collisions

- Ajout de  $e$ , tel que  $h(\text{la\_clé}(e)) = i$  et  $T[i]$  occupé
- Où?
  - à une autre place libre du tableau  $T$
  - dans une zone à part, zone de débordement
- Comment le relier à  $i$ ?
  - par pointeur  $\Rightarrow$  liste chaînée des collisions
  - par calcul  $\Rightarrow$  algorithme déterminant la suite  $\{i_p\}$



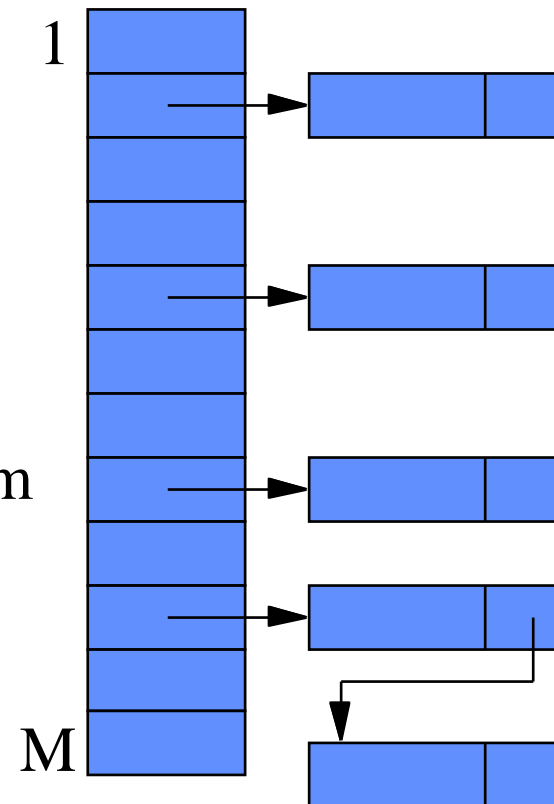
*présentation de la méthode de chaînage en zone de débordement*

simple, efficace et adaptée au disque

# Exemple d'implantation

```
type Place; type Pt_Place is access Place;  
type Place is record  
  Contenu : Element;  
  Suivant : Pt_Place := null;  
end record;  
type Ens_Hach is array (1 .. M) of Pt_Place;
```

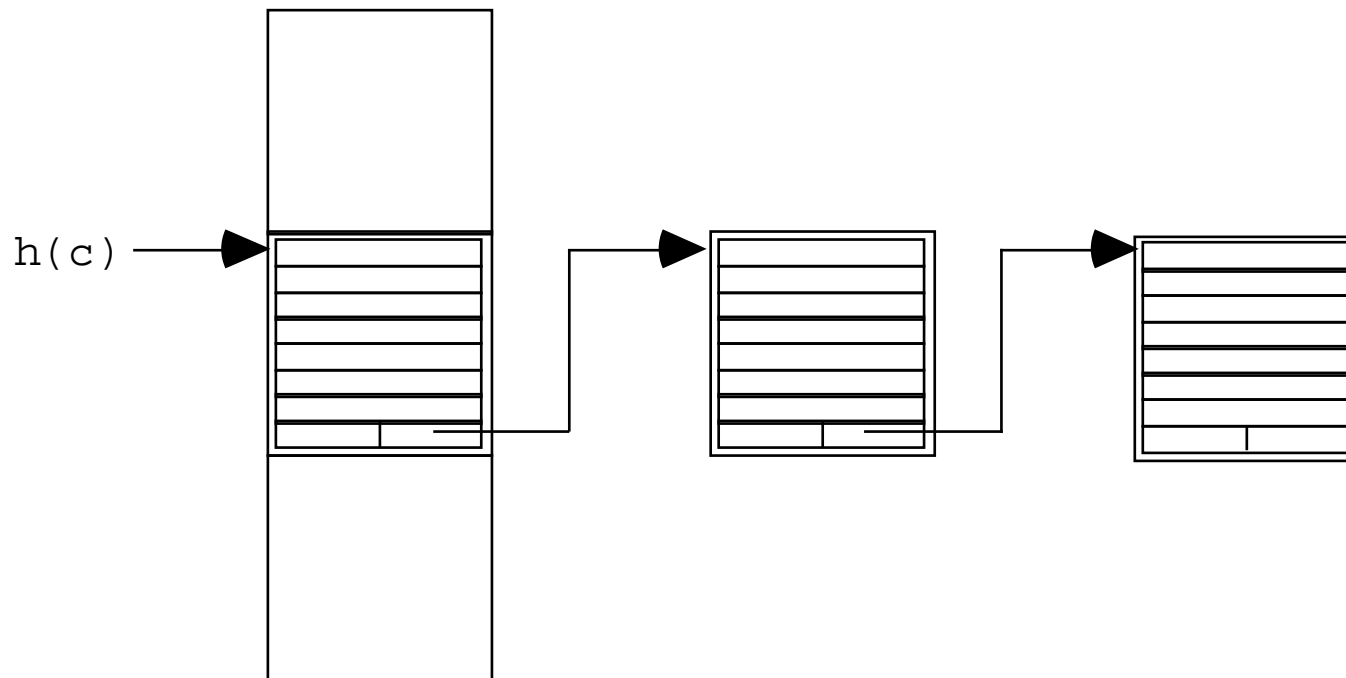
- si fonction hachage correcte, listes sensiblement de même longueur:  $n/m$
- $m$  grand  
 complexité en temps faible  
 perte de place dans T





# Implantation sur disque

- 1 bloc =  $p$  enregistrements de la même liste
- Chaînage des blocs d'une liste



# Implantation en Ada

```
package Hachage_IO is
  type File_Type is limited private;
private
  package L_E is new Listes_Contigues (Element);
  type T_Bloc is record
    Contenu : L_E.Liste (Taille); Suivant : Natural := 0;
  end record;
  package IO is new Ada.Direct_IO (Element_Type => T_Bloc);
  subtype T_Num_B_Hach is Integer range 1..Espace'pos(Espace'Last)
    - Espace'pos (Espace'First) + 2;
  type File_Type is record
    File : IO.File_Type;
    Bloc: T_Bloc; Num_Bloc : Natural := 0;
    Indice : Positive;
    Hachage_Courant : T_Num_B_Hach := T_Num_B_Hach'Last;
  end record;
end Hachage_IO;
```

*chaînage des blocs sur disque*

*transformation espace de hachage en numéros de blocs*

# Les opérations

- Recherche  $\Rightarrow$  parcours de la liste des blocs
- Adjonction  $\Rightarrow$  parcours jusqu'à la fin (contrôle) et ajout au bout
- Suppression  $\Rightarrow$  recherche et tassement ou marquage
- complexité  $\Rightarrow$  nombre de blocs de la liste  $q$ 
  - au pire une seule liste  $\Rightarrow q = n/p$
  - en moyenne  $q = n/mp + 0,5$
  - espace inoccupé:  $1/2q$
- Bon compromis:  $q = 2,5$ , donne 2,5 accès disques et 20% inoccupé  $\Rightarrow m = n/2p$

# Réorganisation

- Quand  $n$  augmente, la taille des listes augmente et il y a dégradation des performances
- Réorganisation par augmentation de  $m$ , avec recopie du contenu du fichier
- Réorganisation dynamique

fonction  $H$ : clé  $0..N$ , et  $h(c) = H(c) \bmod m$

augmentation:  $m' = 2m \Rightarrow h'(c) = h(c) \text{ ou } h(c) + m$

Il est possible de continuer à utiliser le fichier pendant sa réorganisation

- marquer les listes  $m+1$  à  $2m$
- accès liste  $i$  marquée  $\Rightarrow$  demander éclatement liste  $i-m$
- éclatement liste  $i \Rightarrow$  démarquer liste  $i+m$