

情報構造 第十二回

集合・辞書

今日の予定

- 集合 (Set)
 - 集合の仕様
 - 実現：ビットベクトル
 - 実現：連結リスト
 - 辞書 (Dictionary)
 - 辞書の仕様
 - 実現：配列へのベタ詰め法
 - 実現：ハッシュ法
 - オープンハッシュ法
 - クローズドハッシュ法
- } ここがメイン

集合

集合 (set)

- 要素数は**有限**
- 要素の**並び順には意味はない** $\{1, 2, 3\}$ と $\{3, 2, 1\}$ は同じ
- **同じ要素はひとつ**しか表れない
 - $\{1, 2, 2, 3, 1\}$ は集合ではなく, **多重集合** (BAG, multiset)
- 要素間に**線形順序** (全順序) \ll を考えることもある
- 集合S上の線形順序 \ll はつぎのとおりである
 - $\forall a, b \in S$ に対して, $a \ll b$ または $a = b$ または $b \ll a$ の**どれか1つ**だけ成立
 - **排中律**
 - $\forall a, b, c \in S$ に対して, $a \ll b$ かつ $b \ll c$ ならば $a \ll c$ が成立
 - **推移律**

集合の仕様

- **要素：**

- 集合のすべての要素は**同じ型**Elementをもつ
- 要素型Elementが線形順序をもちうる

- **構造：**

- 集合は同じ型の集まり
- 要素の**重複はない**
- 要素を並べる**順序に意味はない**
- 要素数は**有限**で、要素がないものを**空集合**と呼ぶ

- **操作：**

集合に対して

- 和・積・差の集合演算 (Union, Intersection, Difference)
- 要素の挿入・削除・所属 (Insert, Delete, Member)
- 最小要素の取り出し (Min)

など

集合の実現：ビットベクトル

- 表現

- 整数Nの値が小さいとき, $0, 1, \dots, N$ の整数を要素とする集合
- すなわち, 普遍集合 (全体集合) $\{0, 1, 2, \dots, N\}$ 上の集合

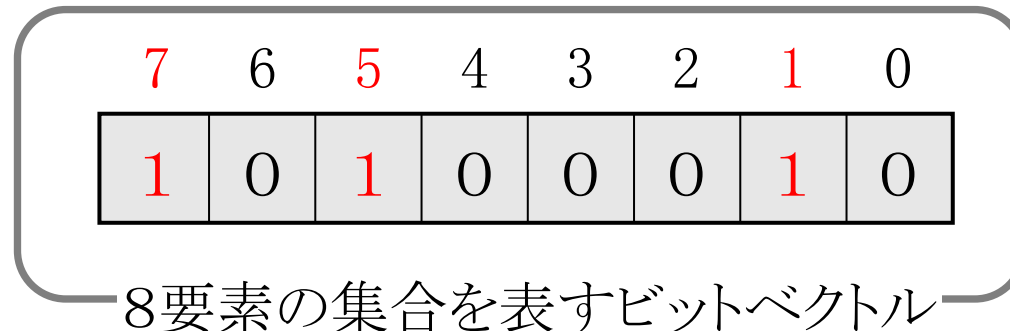
⇒ビットベクトルで実現

- ビットベクトル

- 表現する語・バイトの**ビット上の位置**とその**値0**または**1**で集合の要素を表す
- 例えば, **5**が**集合要素**のとき**ビット位置5**の値が**1**

```
typedef char SET; /* 普遍集合{0, 1, 2, 3, 4, 5, 6, 7} */
```

集合変数
 $S = \{1, 5, 7\}$



実現アルゴリズム：和集合，積集合，差集合

- 集合操作は以下で実現できる

- ビット論理演算子： $|$, $\&$, \sim (or, and, not)
- シフト演算子： $>>$, $<<$ (右シフト，左シフト)

数値の各ビットを右
または左へシフトさ
せる演算子

最下位ビット
least significant
bit (右端ビット)

集合 A の第 i 番目のビットを $A[i]$ で表す (LSB を $i=0$ とする)

- A, B の **和集合** $C = \text{Union}(A, B)$: $C[i] = A[i] | B[i]$
- A, B の **積集合** $C = \text{Intersection}(A, B)$: $C[i] = A[i] \& B[i]$
- A, B の **差集合** $C = \text{Difference}(A, B)$: $C[i] = A[i] \& (\sim B[i])$

| | | | | | | | | |
|---------------|---|---|---|---|---|---|---|---|
| | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
| A = {1, 5, 7} | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

| | | | | | | | | |
|---------------|---|---|---|---|---|---|---|---|
| B = {1, 2, 5} | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
|---------------|---|---|---|---|---|---|---|---|

| | | | | | | | | |
|------------------------|---|---|---|---|---|---|---|---|
| 和集合 $C = A B$ | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
|------------------------|---|---|---|---|---|---|---|---|

| | | | | | | | | |
|-------------------------|---|---|---|---|---|---|---|---|
| 積集合 $C = A \& B$ | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
|-------------------------|---|---|---|---|---|---|---|---|

実現アルゴリズム: 所属, 挿入, 削除

- **所属** : $\text{Member}(i, A) = i \in A$
 - $\{i\}$ を表す整数 2^i と A の **ビット積** ($2^i \& A$) を結果値とする
結果が整数 0 でないとき, すなわち $2^i \& A \neq 0$ のとき $i \in A$
- **挿入** : $\text{Insert}(i, A)$
 - $\{i\}$ を表す整数 2^i と A の **ビット和** ($2^i \mid A$) を結果値とする
- **削除** $\text{Delete}(i, A)$
 - $\{i\}$ を表す整数 2^i の **1の補数** $\sim 2^i$ と A の **ビット積** ($\sim 2^i \& A$) を結果値とする

Memberの例

$A = \{1, 5, 7\}$

$\{5\} = 2^i$

$5 \in A? \quad A \& \{5\}$

| 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |

実現アルゴリズム（ビットベクトル）の効率

- ビットベクトルが1語に納まるとき（ N ビット \leq 1語）
 - **Member, Insert, Delete, Union, Intersection, Difference**すべて—**定時間**
 - …演算レジスタで**数回**の**論理演算**で実行できる
- N ビットに **m** 語要するときは， **m** 倍の**論理演算**が必要となる
- **Min(A)**は， A の**ビット数** N の時間がかかる

普遍集合が列挙型

- 普遍集合を**列挙型**（スカラー型）にすることで，集合要素を**列挙値の識別子**で表せ，集合表現が分かりやすくなる
- 例えば，普遍集合 $\{0, 1, 2, 3, 4, 5, 6, 7\}$ の代わりに
enum color {Black, White, Red, Green, Yellow, Blue, Brown, Purple}
とした場合，列挙値Black, ..., Purpleは，0, ..., 7なので
要素挿入 Member(2, A)は，Member(Red, A)と書ける

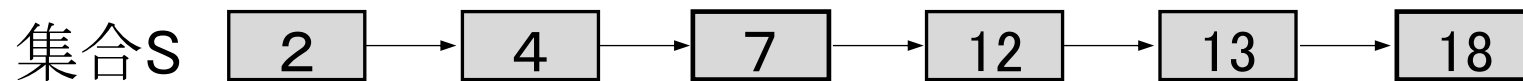
集合値S={Black, Red, Purple}

| Purple | Brown | Blue | Yellow | Green | Red | White | Black |
|--------|-------|------|--------|-------|-----|-------|-------|
| =7 | =6 | =5 | =4 | =3 | =2 | =1 | =0 |
| 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |

集合の実現：連結リスト

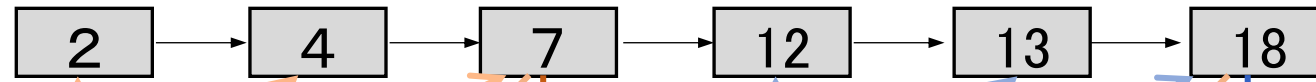
- 表現

- ひとつの**連結リスト**でひとつの**集合**を表す
- 連結リスト中の各構造体で、集合の各要素を表す
- 操作の効率を考えて、集合の要素間の**線形順序**を用い、連結リスト中の**構造体**を**ソート**しておく

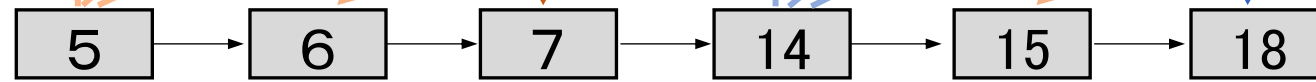


連結リストによる集合演算

集合S



集合T

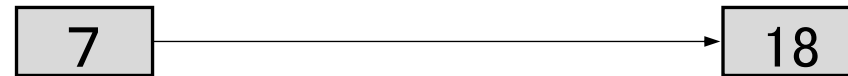


比較

一致

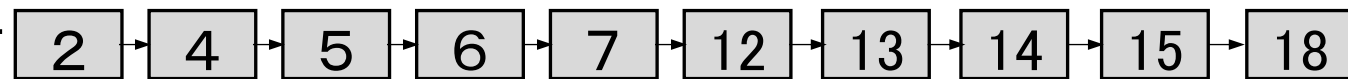
一致した要素を出力

積集合 $S \cap T$



以下の要素を出力

和集合 $S \cup T$



実現アルゴリズム（連結リスト）の効率

- Union, Intersection Differenceの効率
 - 連結リストの構造体が**ソートされている**場合
 - …二つのリストの、**長さの和**
 - 連結リストの構造体が**ソートされていない**場合
 - …二つのリストの、**長さの積**

辭書

辞書 (Dictionary)

- 操作 **Insert, Delete, Member**を備えた**集合**を**辞書**とよぶ

=> 操作Union, Intersection, Differenceを除いた抽象データ型

辞書の仕様

- 要素・構造（辞書の要素と構造は**集合と同じ**）

- 集合のすべての要素は**同じ型**Elementをもつ
- Elementは**線形順序**を持ち得る
- 要素の**重複はない**
- 要素を並べる**順序に意味はない**
- 要素数は**有限**で、要素がないものを**空の辞書**と呼ぶ

例： $\{1,2,3\} = \{3,2,1\} \neq \{3,3,2,1\}$

- 操作（辞書の操作は抽象データ型**集合の部分集合**）

- 要素の**挿入**，**削除**，**所属**（探索）（Insert, Delete, Member）
- など

辞書の操作

辞書型 Dictionary

要素型 Element

辞書型変数 D

要素型データ x

- int **Member**(Element x, Dictionary D)
 - Post: 関数値は $x \in D$ ならば真 (1) さもないと偽 (0)
- **Dictionary Insert**(Element x **Dictionary D**)
 - Post: 関数値は $D \cup \{x\}$
- **Dictionary Delete**(Element x **Dictionary D**)
 - Post: 関数値は $D - \{x\}$
- **Dictionary Create** (void)
 - Post: 関数値は, 空の辞書
- Element **Min**(**Dictionary D**)
 - Pre: $D \neq$ 空の辞書
 - Post: 関数値は D 中の最小要素の値

辞書の実現

- 辞書の**配列**上での実現方法
 - ベタ詰め法
 - ハッシュ法

表現：辞書の実現 配列へのベタ詰め

```
#define maxsize 1000
```

```
/*辞書の要素の最大数*/
```

```
typedef struct{
```

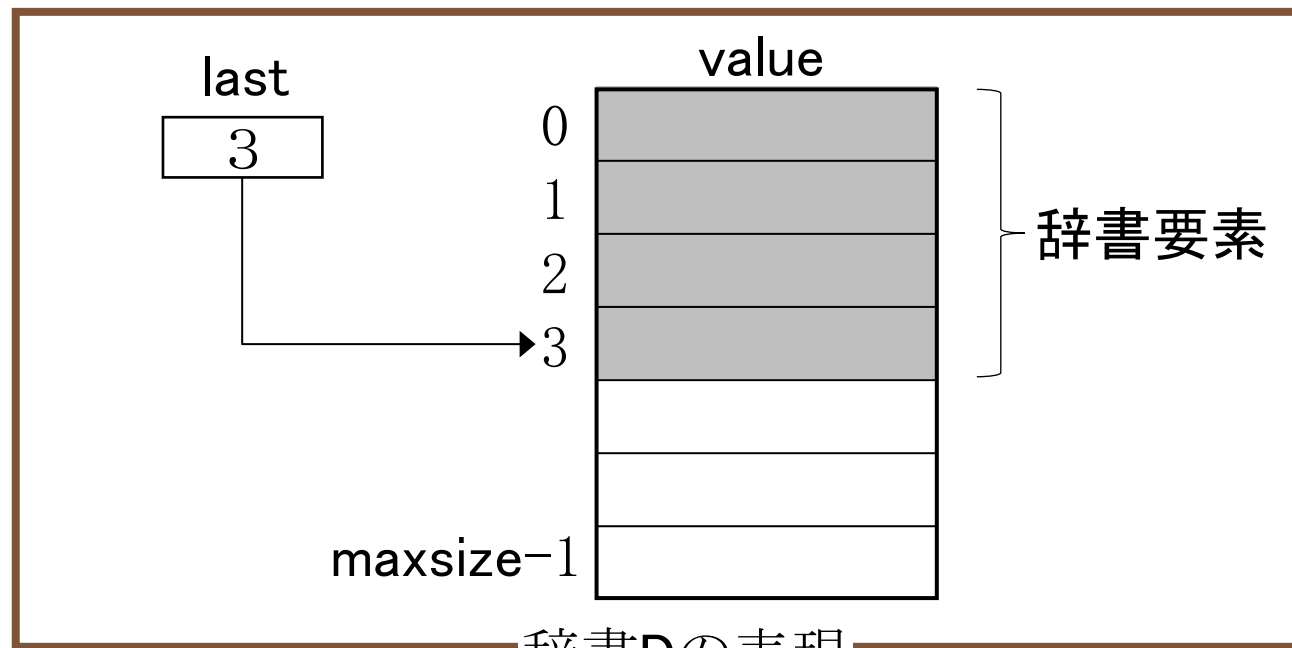
```
    Element value[maxsize];
```

```
    int last;
```

```
/*辞書要素の最後, -1の時は空*/
```

```
}Dictionary;
```

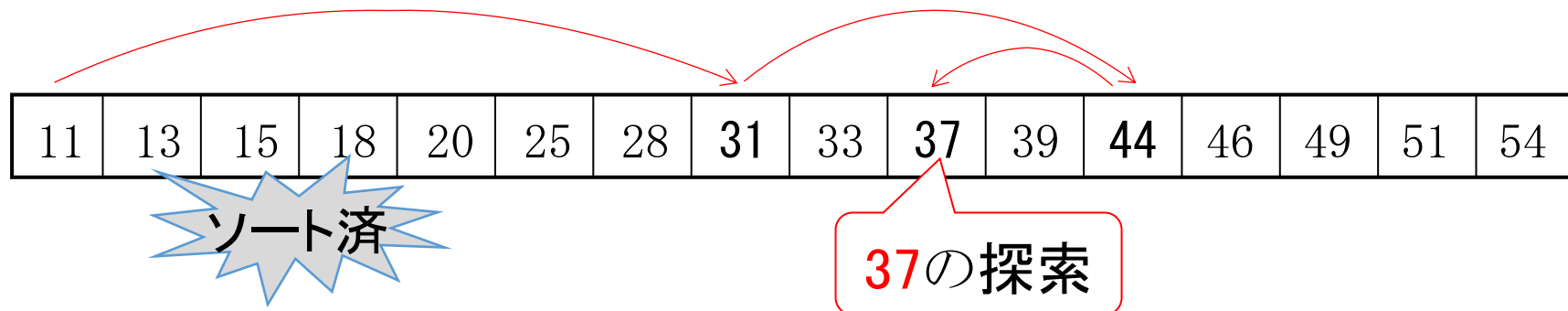
0からlastが
辞書要素



辞書Dの表現

効率：辞書の実現 配列へのベタ詰め

- 辞書の要素を**ソートしない**とき
 - Insert, Delete, Memberは, **最悪 N** (辞書の要素数) かかる
- 辞書の要素を**ソートした**とき
 - Member, Insert, Deleteは, 探索に最悪 $\log_2 N$ かかる



辞書の実現：ハッシュ法

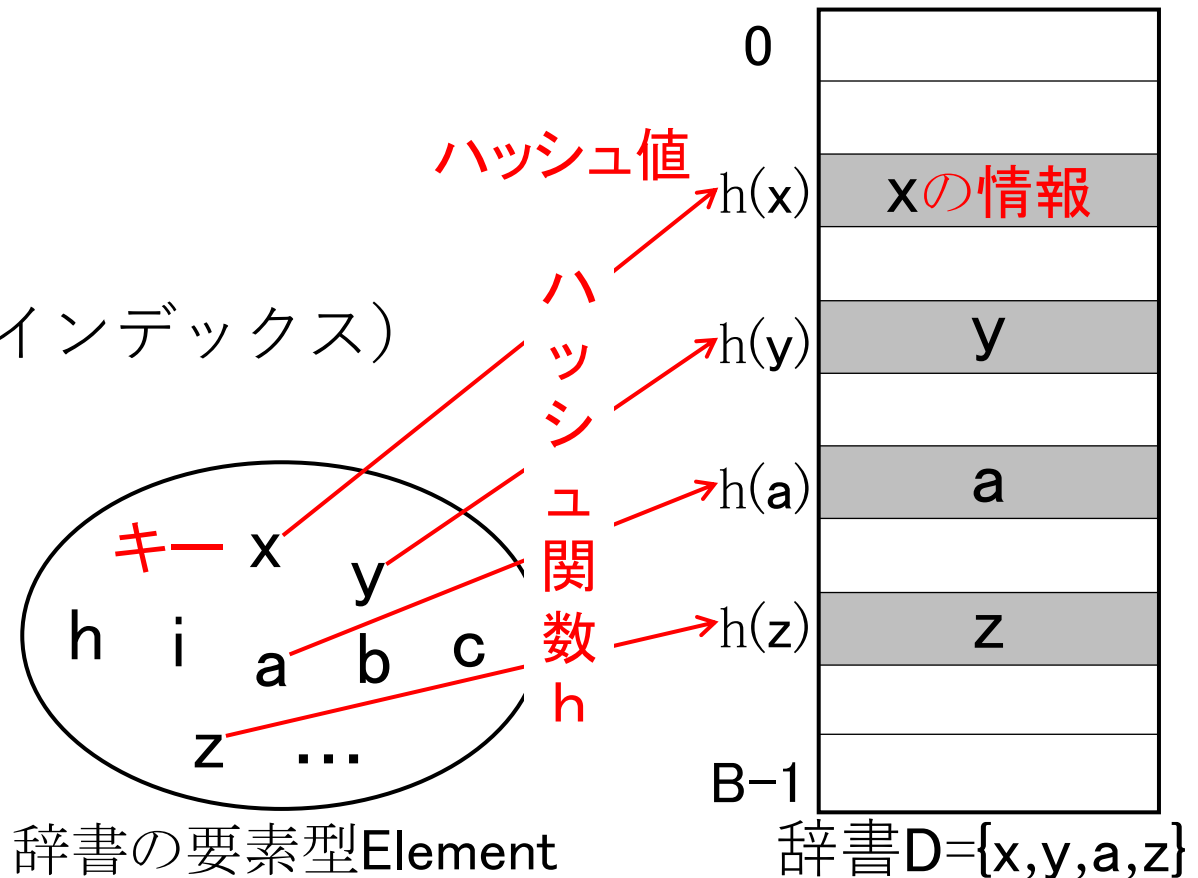
- 辞書の**配列へのベタ詰め**の実現は
 - 操作は**最悪**で**登録要素数 N** （ソートすれば **$\log_2 N$** ）のオーダーの時間計算量
- ハッシュ法
 - 辞書を配列（ハッシュ表）で実現
 - 各操作の**平均の時間計算量を一定**にする

辞書要素とハッシュ表の対応

- 辞書の要素型Element
- 辞書要素を格納する配列T（ハッシュ表）…辞書の表現 **ハッシュ表 T**
- Tのインデックスが0からB-1

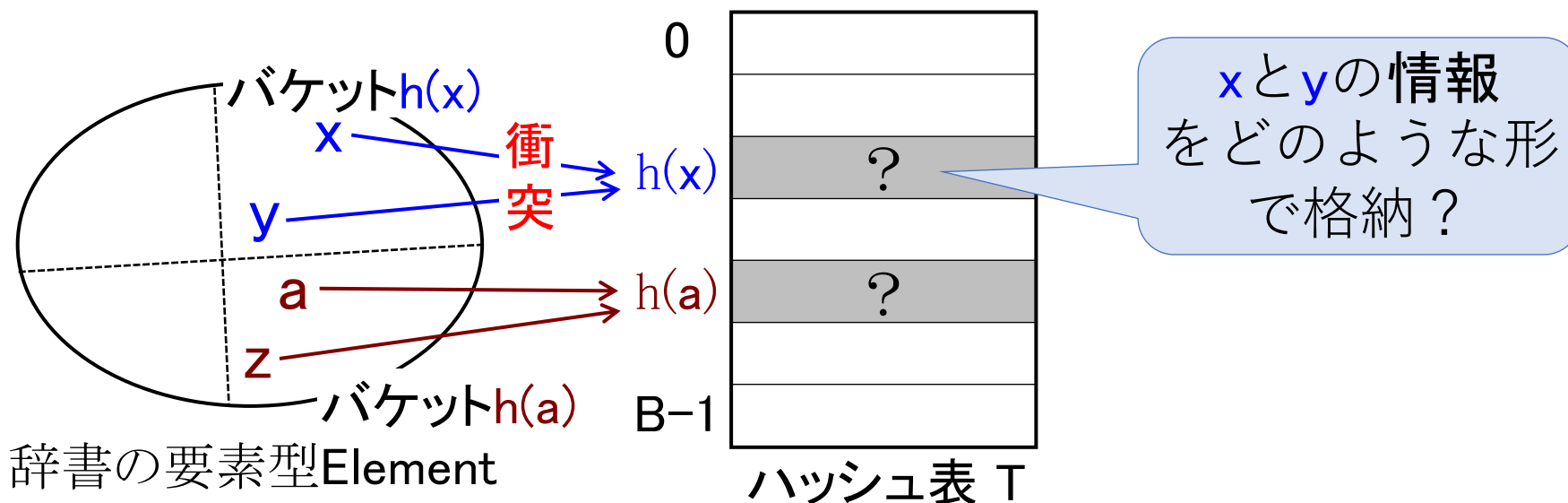
- 次の関数hを考える

- $h: \text{Element} \rightarrow [0 \cdots B-1]$ （要素からインデックス）
- hを**ハッシュ関数**と呼ぶ
- Element型の要素xをキー
- $h(x)$ をキーxの**ハッシュ値**と呼ぶ



衝突 (collision)

- Elementすべての要素のハッシュ値が異なる (h が単斜) なら辞書要素 (キー) x の情報は配列 T の第 $h(x)$ 番目に格納できる?
- ハッシュ表を非常に大きくしなければ, 一意の関数を見つけるのは難しい!



- いくつかの要素が**同じハッシュ値**を持つ
 - このクラスを**バケット**, そのハッシュ値を**バケット番号**という
- 同じバケット内の異なる辞書要素 (キー) x と y は, **衝突**するという

ハッシュ法の種類

- 要素の衝突に多する処理はいろいろあるが、ここでは二つの方法を取り上げる

1. オープンハッシュ法 (チェーン法)

2. クローズドハッシュ法 (オープンアドレス法)

オープンハッシュ法

(open hashing)

overflow hash

direct chaining

chaining

外部ハッシュ

⇔

クローズドハッシュ法

(closed hashing)

⇔

open hash

⇔

open addressing

⇔

open addressing

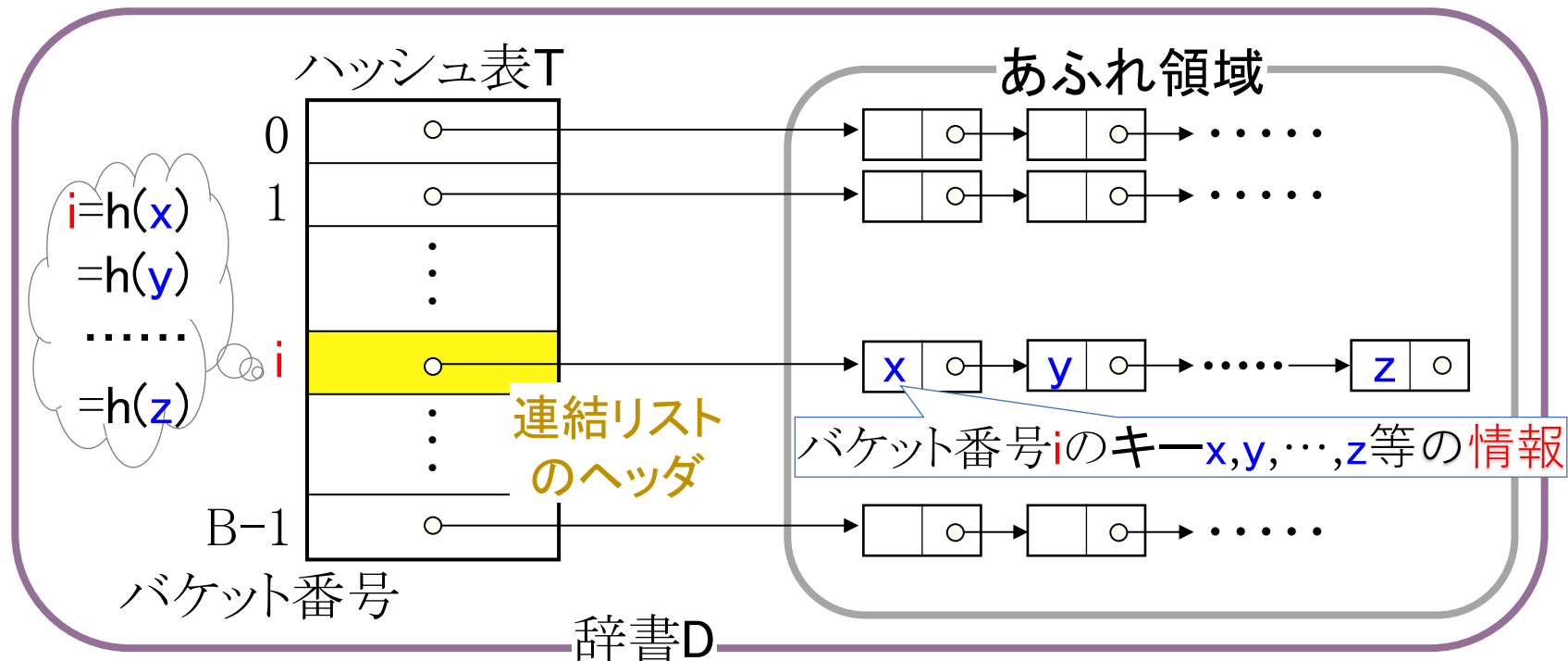
⇔

内部ハッシュ

オープンが逆になる
他の名前があるので
注意

オープンハッシュ法：辞書の表現

- 配列T： ハッシュ表
インデックスがバケット番号 $0, 1, \dots, B-1$
- 同じバケット i 中の衝突要素（衝突キー）：
T[i]をヘッダとする**連結リスト**に登録（Tの**要素**はキーではない）

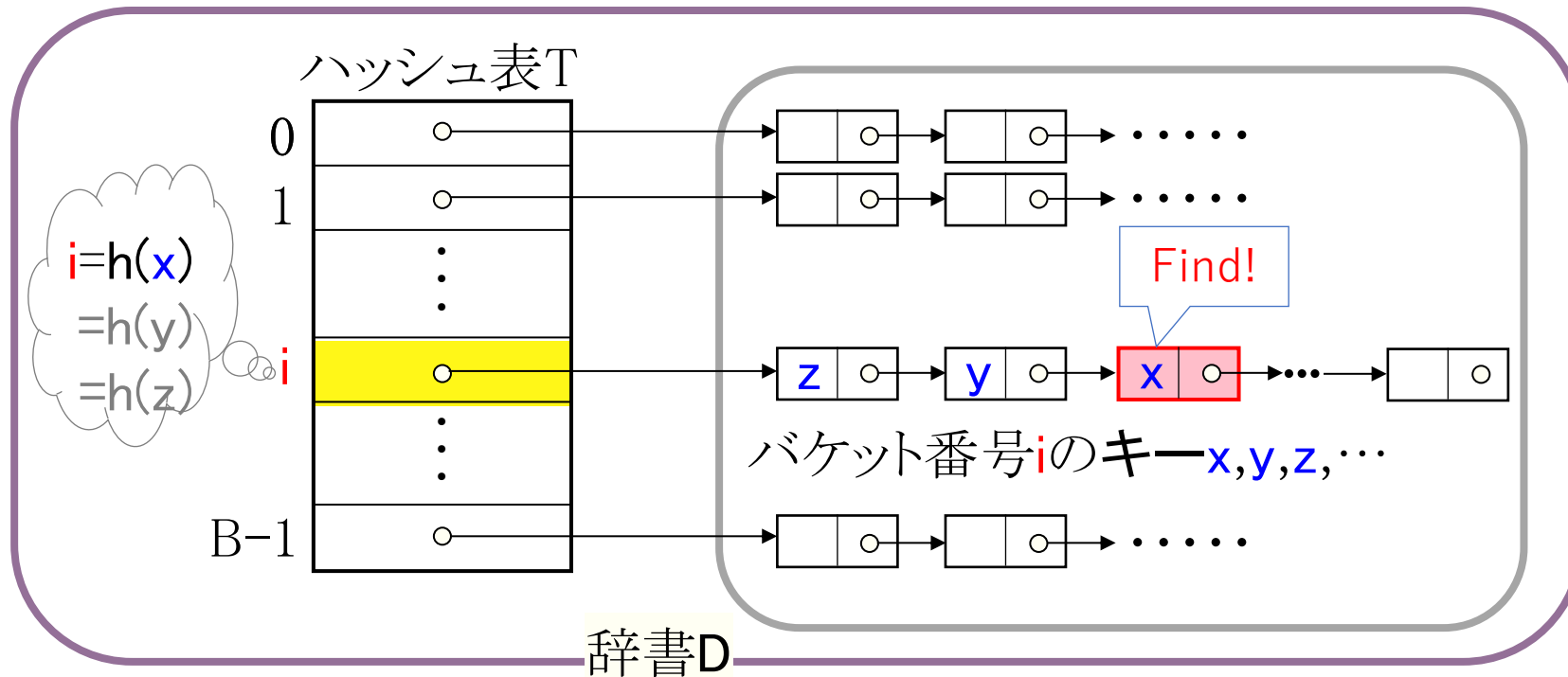


オープンハッシュ法：操作

- 挿入：Insert(x , D)
- 削除：Delete(x , D)
- 所属・探索：Member(x , D)

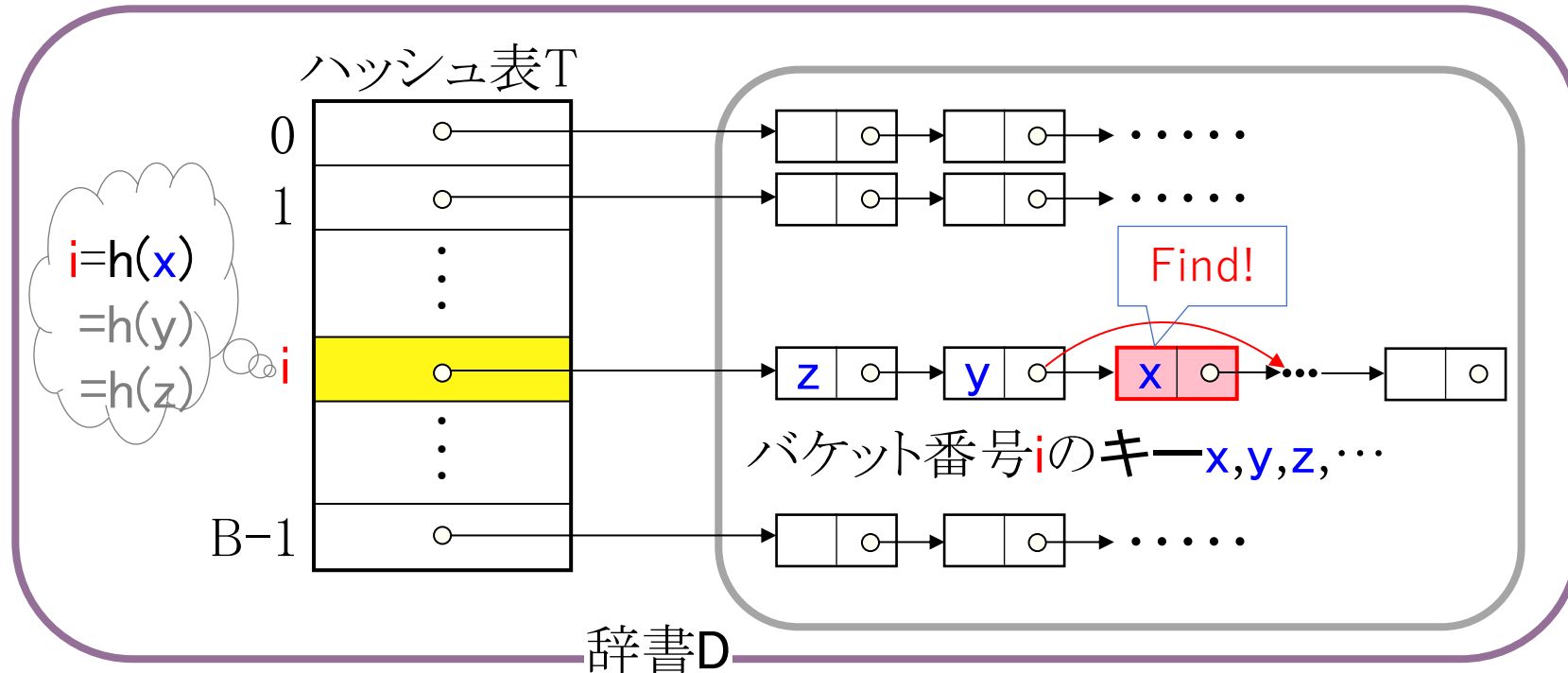
実現アルゴリズム: 挿入 Insert (x, D)

1. キー（辞書要素） x のバケット i ($=h(x)$)を求める
2. ハッシュ表の要素 $T[i]$ をヘッダとする連結リスト中に x があるか調べる
3. あれば, 何もせず終了
4. なければ, x を連結リストの要素として挿入して終了



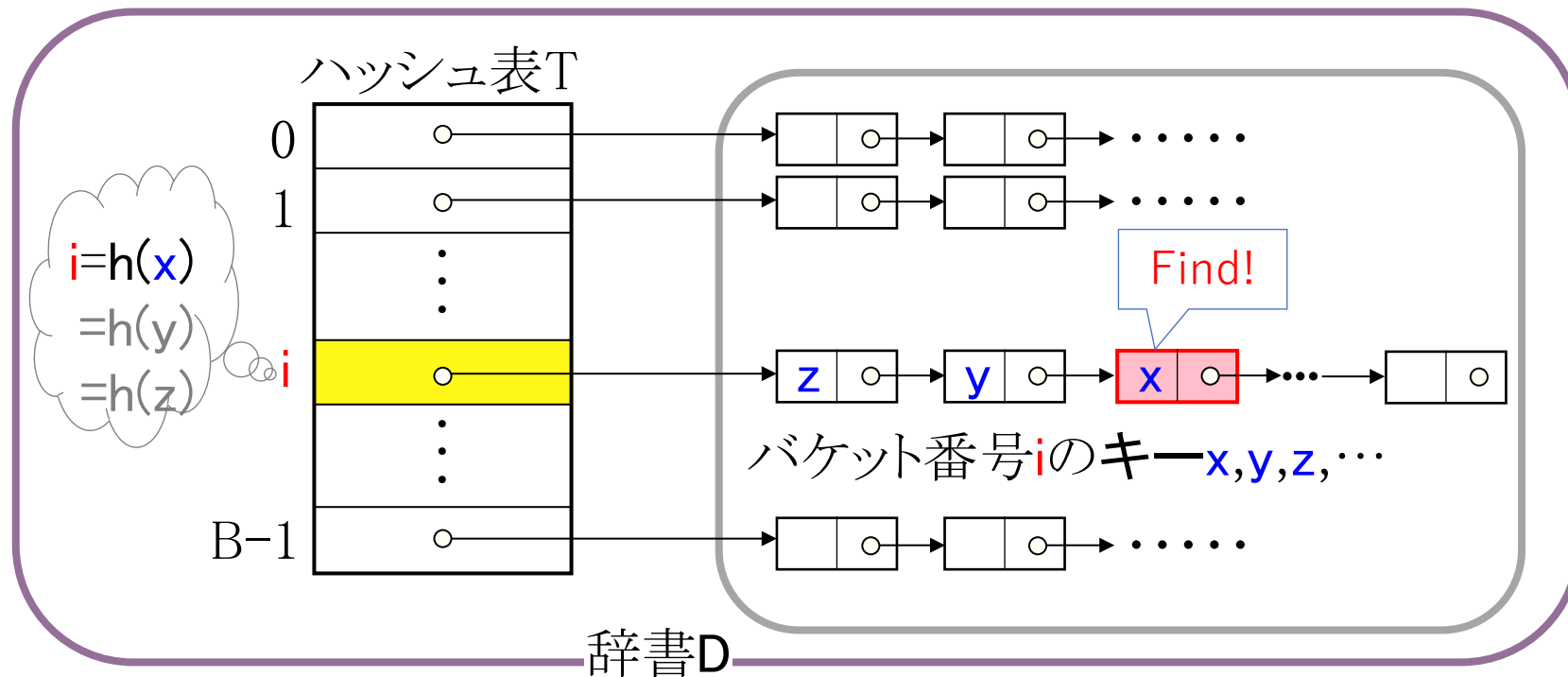
実現アルゴリズム: 削除Delete (x, D)

1. キー x のバケット i ($=h(x)$)を求める
2. ハッシュ表の要素 $T[i]$ をヘッダとする連結リスト中に x があるか調べる
3. あれば, 連結リストから削除して終了
4. なければ, 何もせず終了



実現アルゴリズム: 探索Member (x, D)

1. キー x のバケット i ($=h(x)$) を求める
2. ハッシュ表の要素 $T[i]$ をヘッダとする連結リスト中に x があるか調べ, その結果を関数値として終了



ハッシュ関数の選定

- **キー**（辞書要素）を各バケットに**均一**に**割り当てる**ハッシュ関数が好ましい！
- **キー****x**が**文字列**の時は，**数値**に変換してから，バケットに割り当てる

ハッシュ関数の例

- キー x : 文字列 $c_{n-1} \cdots c_1 c_0$
- バケット数 B : 50
- ハッシュ値: 「各文字の**ASCIIコードの和**の B による**剰余**」

$$h(c_{n-1} \cdots c_1 c_0) = (c_{n-1} + \cdots + c_1 + c_0) \% B$$

$$h(\text{"A00"}) = (65 + 48 + 48) \% 50 = 11$$

C言語では、文字の和は、
文字のASCIIコードの和

- キー100個: A00, A01, ..., A99を
50個のバケットに割り当てる
=> キーが**均一**に割り当てられない

- キー集合{"A00", ..., "A99"}, $h(x)$ = 「文字コードの和のBの剰余」

| | | |
|--------|-----|---|
| バケット11 | 1個 | A00 |
| バケット12 | 2個 | A01 A10 |
| バケット13 | 3個 | A02 A11 A20 |
| バケット14 | 4個 | A03 A12 A21 A30 |
| バケット15 | 5個 | A04 A13 A22 A31 A40 |
| バケット16 | 6個 | A05 A14 A23 A32 A41 A50 |
| バケット17 | 7個 | A06 A15 A24 A33 A42 A51 A60 |
| バケット18 | 8個 | A07 A16 A25 A34 A43 A52 A61 A70 |
| バケット19 | 9個 | A08 A17 A26 A35 A44 A53 A62 A71 A80 |
| バケット20 | 10個 | A09 A18 A27 A36 A45 A54 A63 A72 A81 A90 |
| バケット21 | 9個 | A19 A28 A37 A46 A55 A64 A73 A82 A91 |
| バケット22 | 8個 | A29 A38 A47 A56 A65 A74 A83 A92 |
| バケット23 | 7個 | A39 A48 A57 A66 A75 A84 A93 |
| バケット24 | 6個 | A49 A58 A67 A76 A85 A94 |
| バケット25 | 5個 | A59 A68 A77 A86 A95 |
| バケット26 | 4個 | A69 A78 A87 A96 |
| バケット27 | 3個 | A79 A88 A97 |
| バケット28 | 2個 | A89 A98 |
| バケット29 | 1個 | A99 |

バケット0から10
のキーは0個

バケット30から49
のキーは0個

平均 2
分散 9.4

ハッシュ関数の例

- キー x : 文字列 $c_{n-1} \cdots c_1 c_0$
- バケット数 B : 50
- ハッシュ値: 「 n 個の文字列をASCIIコードで128進数 n 桁の数値」の
 B による剰余

$$h(c_{n-1} \cdots c_1 c_0) = (c_{n-1} * 128^{n-1} + \cdots + c_1 * 128^1 + c_0 * 128^0) \% B$$

$$\begin{aligned} \text{例: } h(A49) &= ('A' * 128^2 + '4' * 128^1 + '9' * 128^0) \% B \\ &= (65 * 16381 + 52 * 128 + 57 * 1) \% B \\ &= 1071673 \% 50 \\ &= 23 \end{aligned}$$

- キー100個: A00, A01, ..., A99を 50個のバケットに割り当てる
=> キーがほぼ**均一**に**割り当てられる**

128進数3桁ならば, $128^3 (= 2,097,152)$ の大きさのハッシュ表を用意すれば**ハッシュ値は単斜** (衝突は起こらない)

- キー集合{“A00”, ..., “A99”}, $h(x)$ = 「128進数のBの剰余」
 - すべてのバケットで、キーは1個から3個

| | | |
|--------|----|-------------|
| バケット0 | 2個 | A58 A72 |
| バケット1 | 2個 | A59 A73 |
| バケット2 | 2個 | A00 A74 |
| バケット3 | 2個 | A01 A75 |
| バケット4 | 3個 | A02 A76 A90 |
| バケット5 | 3個 | A03 A77 A91 |
| バケット6 | 3個 | A04 A78 A92 |
| バケット7 | 3個 | A05 A79 A93 |
| バケット8 | 3個 | A06 A20 A94 |
| バケット9 | 3個 | A07 A21 A95 |
| バケット10 | 3個 | A08 A22 A96 |
| バケット11 | 3個 | A09 A23 A97 |
| バケット12 | 2個 | A24 A98 |
| バケット13 | 2個 | A25 A99 |
| バケット14 | 2個 | A26 A40 |
| バケット15 | 2個 | A27 A41 |
| バケット16 | 2個 | A28 A42 |
| バケット17 | 2個 | A29 A43 |
| バケット18 | 1個 | A44 |
| バケット19 | 1個 | A45 |
| バケット20 | 2個 | A46 A60 |
| バケット21 | 2個 | A47 A61 |
| バケット22 | 2個 | A48 A62 |
| バケット23 | 2個 | A49 A63 |
| バケット24 | 1個 | A64 |

| | | |
|--------|----|---------|
| バケット25 | 1個 | A65 |
| バケット26 | 2個 | A66 A8 |
| バケット27 | 2個 | A67 A81 |
| バケット28 | 2個 | A68 A82 |
| バケット29 | 2個 | A69 A83 |
| バケット30 | 2個 | A10 A84 |
| バケット31 | 2個 | A11 A85 |
| バケット32 | 2個 | A12 A86 |
| バケット33 | 2個 | A13 A87 |
| バケット34 | 2個 | A14 A88 |
| バケット35 | 2個 | A15 A89 |
| バケット36 | 2個 | A16 A30 |
| バケット37 | 2個 | A17 A31 |
| バケット38 | 2個 | A18 A32 |
| バケット39 | 2個 | A19 A33 |
| バケット40 | 1個 | A34 |
| バケット41 | 1個 | A35 |
| バケット42 | 2個 | A36 A50 |
| バケット43 | 2個 | A37 A51 |
| バケット44 | 2個 | A38 A52 |
| バケット45 | 2個 | A39 A53 |
| バケット46 | 1個 | A54 |
| バケット47 | 1個 | A55 |
| バケット48 | 2個 | A56 A70 |
| バケット49 | 2個 | A57 A71 |

| | |
|----|------|
| 平均 | 2 |
| 分散 | 0.32 |

ハッシュ関数の作り方

- **除算法** $h(x) = \text{キー}x \text{を数値化} \% B$

- **平方採中法（乗算法，中央2乗法）**

$$h(x) = (n^2 / C) \% B$$

擬似乱数の生成方法

- 文字キー x を変換した数値 n が0から K とする
- $B \cdot C^2 \doteq K^2$ となるように整数 C で，数値 n の2乗の中央値をとることができる
- 原理は，次のようにキーの数値54321の2乗の中央値077をとることである。
これは，中央の桁が全桁に依存することによる

$$54321^2 = 2950771041$$

- 折り込み法

$$\underline{12345678} = 0001 | 2345 | 678$$

$$\Rightarrow 1000 + 2345 + 867 = 4221$$

• 平方採中法 キー集合 $\{0, \dots, 499\}$, $h(x) = (n^2 / C) \% B$, $K = 499$, $C = 71$

| | | | | | |
|---------|-----|---------|-----|---------|-----|
| バケツト0: | 14個 | バケツト17: | 9個 | バケツト34: | 10個 |
| バケツト1: | 16個 | バケツト18: | 8個 | バケツト35: | 8個 |
| バケツト2: | 17個 | バケツト19: | 12個 | バケツト36: | 9個 |
| バケツト3: | 8個 | バケツト20: | 7個 | バケツト37: | 12個 |
| バケツト4: | 12個 | バケツト21: | 12個 | バケツト38: | 11個 |
| バケツト5: | 9個 | バケツト22: | 5個 | バケツト39: | 8個 |
| バケツト6: | 13個 | バケツト23: | 8個 | バケツト40: | 7個 |
| バケツト7: | 9個 | バケツト24: | 10個 | バケツト41: | 9個 |
| バケツト8: | 11個 | バケツト25: | 10個 | バケツト42: | 14個 |
| バケツト9: | 13個 | バケツト26: | 16個 | バケツト43: | 5個 |
| バケツト10: | 8個 | バケツト27: | 10個 | バケツト44: | 10個 |
| バケツト11: | 11個 | バケツト28: | 8個 | バケツト45: | 12個 |
| バケツト12: | 5個 | バケツト29: | 8個 | バケツト46: | 9個 |
| バケツト13: | 8個 | バケツト30: | 11個 | バケツト47: | 11個 |
| バケツト14: | 13個 | バケツト31: | 10個 | バケツト48: | 3個 |
| バケツト15: | 14個 | バケツト32: | 6個 | バケツト49: | 10個 |
| バケツト16: | 12個 | バケツト33: | 9個 | | |

| | |
|----|-----|
| 平均 | 10 |
| 分散 | 8.7 |

• 平方採中法 キー集合{"A00", ..., "A99"}, $h(x) = (n^2 / C) \% B$

• $n = x[0] * 128^2 + x[1] * 128 + x[2]$, $K = 1072313("A99" - 1)$, $C = 153187$

| | | | | | |
|---------|----|---------|----|---------|----|
| バケツト0: | 1個 | バケツト18: | 2個 | バケツト36: | 1個 |
| バケツト1: | 4個 | バケツト19: | 3個 | バケツト37: | 3個 |
| バケツト2: | 2個 | バケツト20: | 1個 | バケツト38: | 2個 |
| バケツト3: | 0個 | バケツト21: | 4個 | バケツト39: | 1個 |
| バケツト4: | 2個 | バケツト22: | 1個 | バケツト40: | 1個 |
| バケツト5: | 3個 | バケツト23: | 3個 | バケツト41: | 3個 |
| バケツト6: | 1個 | バケツト24: | 2個 | バケツト42: | 1個 |
| バケツト7: | 4個 | バケツト25: | 1個 | バケツト43: | 5個 |
| バケツト8: | 1個 | バケツト26: | 1個 | バケツト44: | 1個 |
| バケツト9: | 2個 | バケツト27: | 4個 | バケツト45: | 1個 |
| バケツト10: | 2個 | バケツト28: | 1個 | バケツト46: | 2個 |
| バケツト11: | 2個 | バケツト29: | 4個 | バケツト47: | 2個 |
| バケツト12: | 0個 | バケツト30: | 1個 | バケツト48: | 0個 |
| バケツト13: | 4個 | バケツト31: | 1個 | バケツト49: | 4個 |
| バケツト14: | 1個 | バケツト32: | 2個 | | |
| バケツト15: | 4個 | バケツト33: | 2個 | | |
| バケツト16: | 1個 | バケツト34: | 0個 | | |
| バケツト17: | 1個 | バケツト35: | 5個 | | |

| | |
|----|-----|
| 平均 | 2 |
| 分散 | 1.8 |

• 平方採中法 キー集合{"A00", ..., "A99"}, $h(x) = (n^2 / C) \% B$

• $n = x[0] * 128^2 + x[1] * 128 + x[2]$, $K = 1161("A99" - "A00")$, $C = 165$

| | | | | | |
|---------|----|---------|----|---------|----|
| バケット0: | 3個 | バケット18: | 1個 | バケット36: | 2個 |
| バケット1: | 1個 | バケット19: | 1個 | バケット37: | 2個 |
| バケット2: | 1個 | バケット20: | 4個 | バケット38: | 3個 |
| バケット3: | 2個 | バケット21: | 3個 | バケット39: | 3個 |
| バケット4: | 0個 | バケット22: | 0個 | バケット40: | 0個 |
| バケット5: | 5個 | バケット23: | 2個 | バケット41: | 1個 |
| バケット6: | 0個 | バケット24: | 3個 | バケット42: | 0個 |
| バケット7: | 0個 | バケット25: | 3個 | バケット43: | 0個 |
| バケット8: | 2個 | バケット26: | 2個 | バケット44: | 1個 |
| バケット9: | 1個 | バケット27: | 3個 | バケット45: | 3個 |
| バケット10: | 2個 | バケット28: | 3個 | バケット46: | 3個 |
| バケット11: | 4個 | バケット29: | 1個 | バケット47: | 3個 |
| バケット12: | 4個 | バケット30: | 2個 | バケット48: | 1個 |
| バケット13: | 2個 | バケット31: | 2個 | バケット49: | 0個 |
| バケット14: | 2個 | バケット32: | 7個 | | |
| バケット15: | 2個 | バケット33: | 0個 | | |
| バケット16: | 4個 | バケット34: | 1個 | | |
| バケット17: | 2個 | バケット35: | 3個 | | |

| | |
|----|------|
| 平均 | 2 |
| 分散 | 2.16 |

• 平方採中法 キー集合{"A00", ..., "A99"}, $h(x) = (n^2 / C) \% B$

• $n = x[0] * 128^2 + x[1] * 128 + x[2]$, $K = 2097157$ (128進数3桁), $C = 299593$

| | | | | | |
|----------|----|----------|----|----------|----|
| バケット 0: | 3個 | バケット 18: | 2個 | バケット 36: | 4個 |
| バケット 1: | 3個 | バケット 19: | 1個 | バケット 37: | 2個 |
| バケット 2: | 2個 | バケット 20: | 2個 | バケット 38: | 0個 |
| バケット 3: | 2個 | バケット 21: | 2個 | バケット 39: | 3個 |
| バケット 4: | 2個 | バケット 22: | 1個 | バケット 40: | 1個 |
| バケット 5: | 1個 | バケット 23: | 3個 | バケット 41: | 3個 |
| バケット 6: | 1個 | バケット 24: | 1個 | バケット 42: | 0個 |
| バケット 7: | 3個 | バケット 25: | 2個 | バケット 43: | 5個 |
| バケット 8: | 2個 | バケット 26: | 3個 | バケット 44: | 2個 |
| バケット 9: | 3個 | バケット 27: | 2個 | バケット 45: | 1個 |
| バケット 10: | 3個 | バケット 28: | 3個 | バケット 46: | 3個 |
| バケット 11: | 2個 | バケット 29: | 1個 | バケット 47: | 2個 |
| バケット 12: | 1個 | バケット 30: | 3個 | バケット 48: | 2個 |
| バケット 13: | 2個 | バケット 31: | 2個 | バケット 49: | 0個 |
| バケット 14: | 2個 | バケット 32: | 1個 | | |
| バケット 15: | 2個 | バケット 33: | 1個 | | |
| バケット 16: | 3個 | バケット 34: | 3個 | | |
| バケット 17: | 2個 | バケット 35: | 0個 | | |

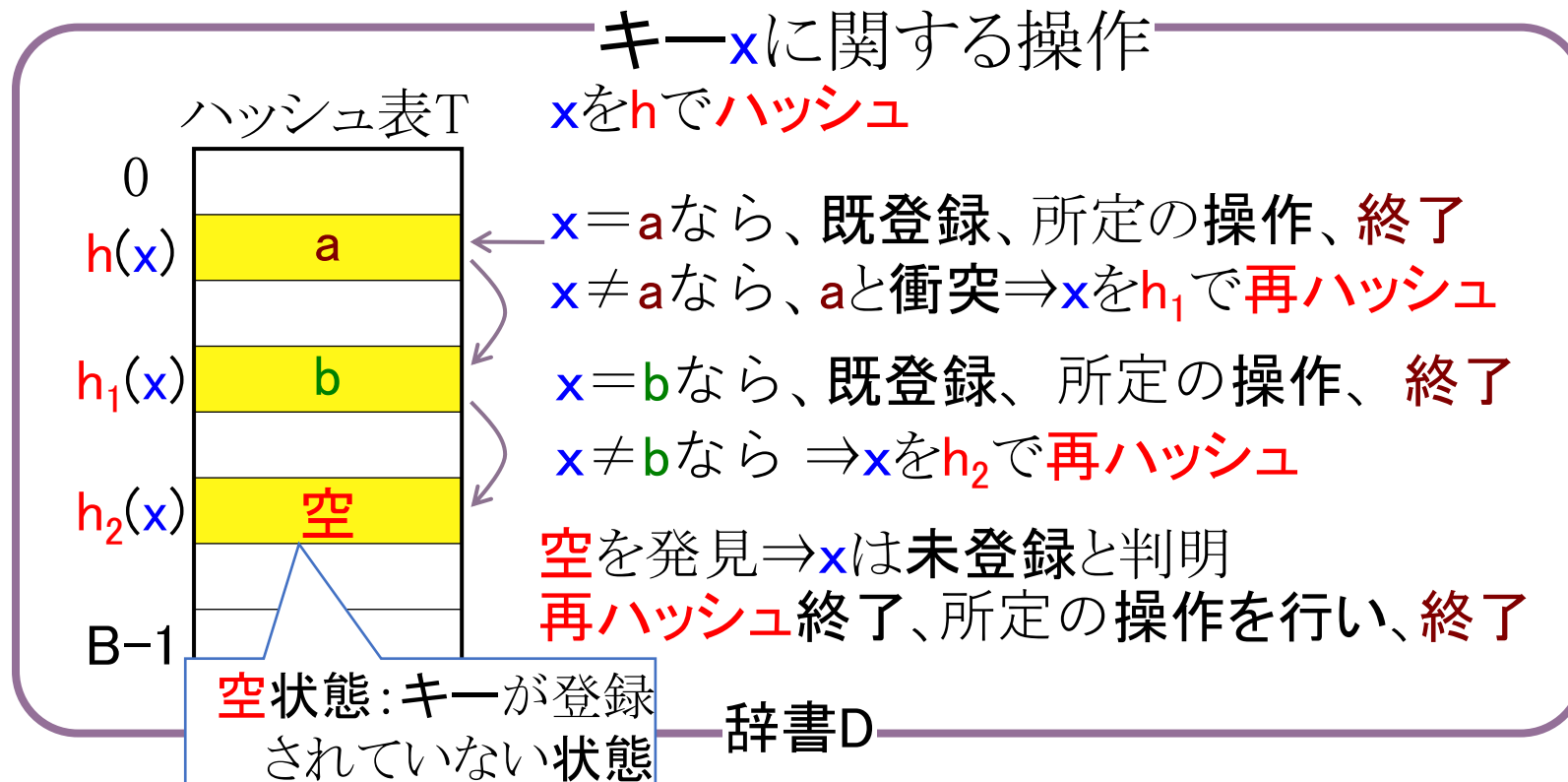
| | |
|----|------|
| 平均 | 2 |
| 分散 | 1.08 |

オープンハッシュの平均時間計算量

- バケット数 B , 登録要素数 N
⇒1バケットあたりの平均要素数は N/B
- **Member, Insert, Delete**の操作1回に要する平均の時間計算量は $O(1 + N/B)$
↓
 $B \div N \Rightarrow$ **平均時間計算量一定**  一つの要素にバケット一つ
- **表の再構成**
 - 登録要素数がバケット数 B の2倍以上になったとき, バケット数が2倍の新しいハッシュ表を作る
⇒いつも $(N/B) < 2$ となるので, 操作の平均効率率は**3未満**
⇒表の大きさ B は登録要素数 N の2倍程度
- **最小の要素**を取り出す操作**Min**の効率のよいアルゴリズムはない

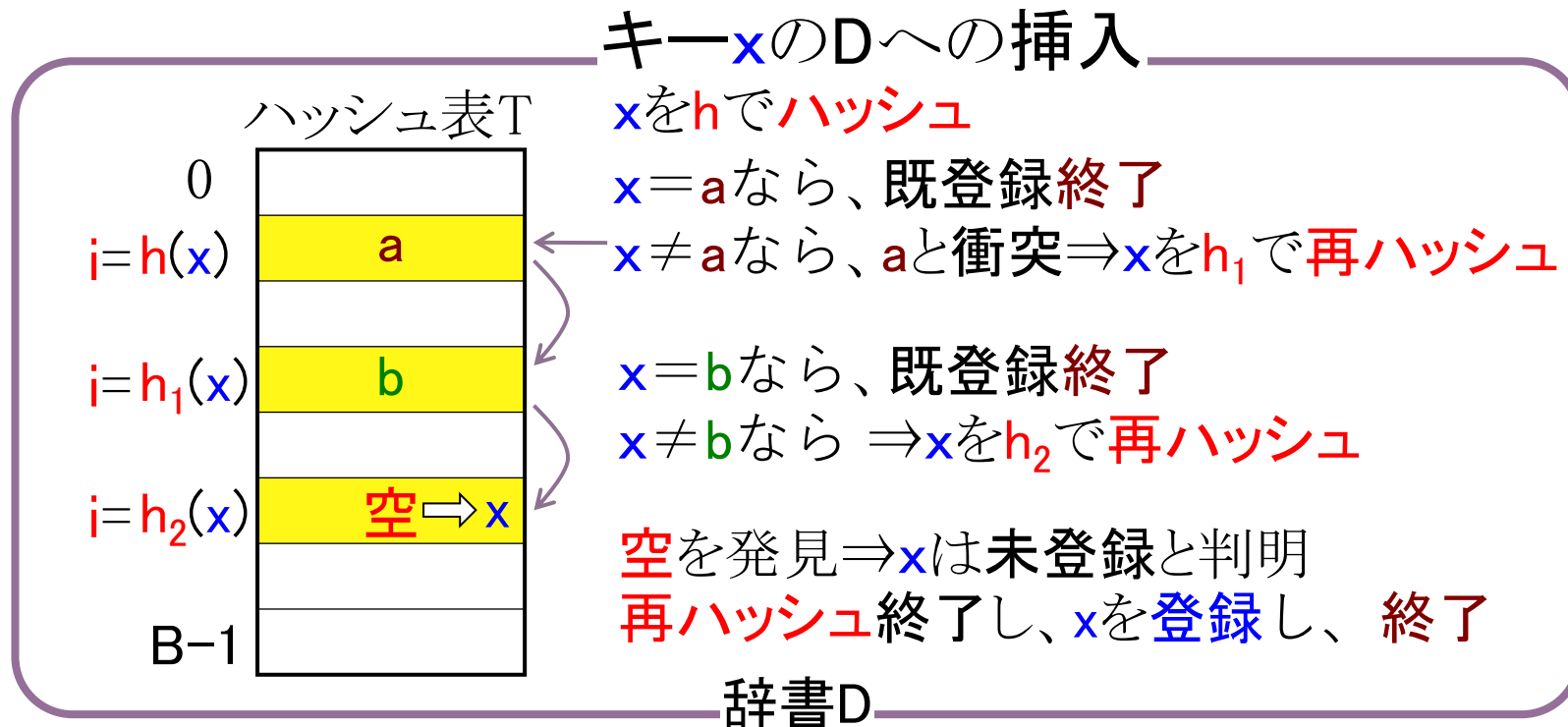
クローズドハッシュ法：辞書の表現

- ハッシュ表Tに, 直接キー (辞書要素) **そのもの**を格納
- ハッシュ表Tに, ハッシュ関数**h**以外に, **再ハッシュ関数** h_1, h_2, h_3, \dots (高々 $B-1$ 個) を使う



削除操作のない場合の挿入Insert(x, D)

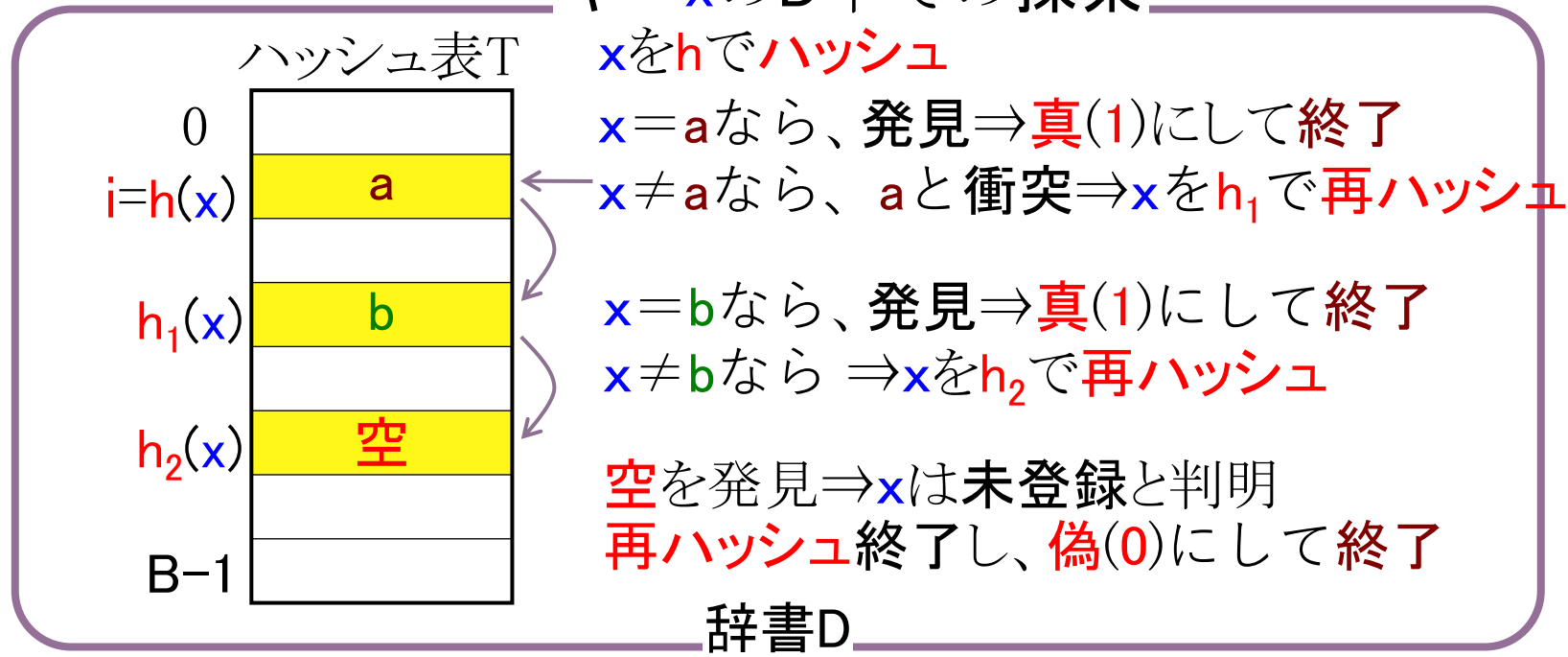
1. キー x をハッシュとするハッシュ値 $i = h(x)$ を求める
2. ハッシュ表の要素 $T[i]$ が
 - x ならば, 既登録より終了
 - 空の状態ならば, 未登録なので, ここに x を登録し終了
 - x 以外のキーならば, 再ハッシュし再ハッシュ値 i を求め2.を繰り返す



削除操作のない場合の所属Member (x, D)

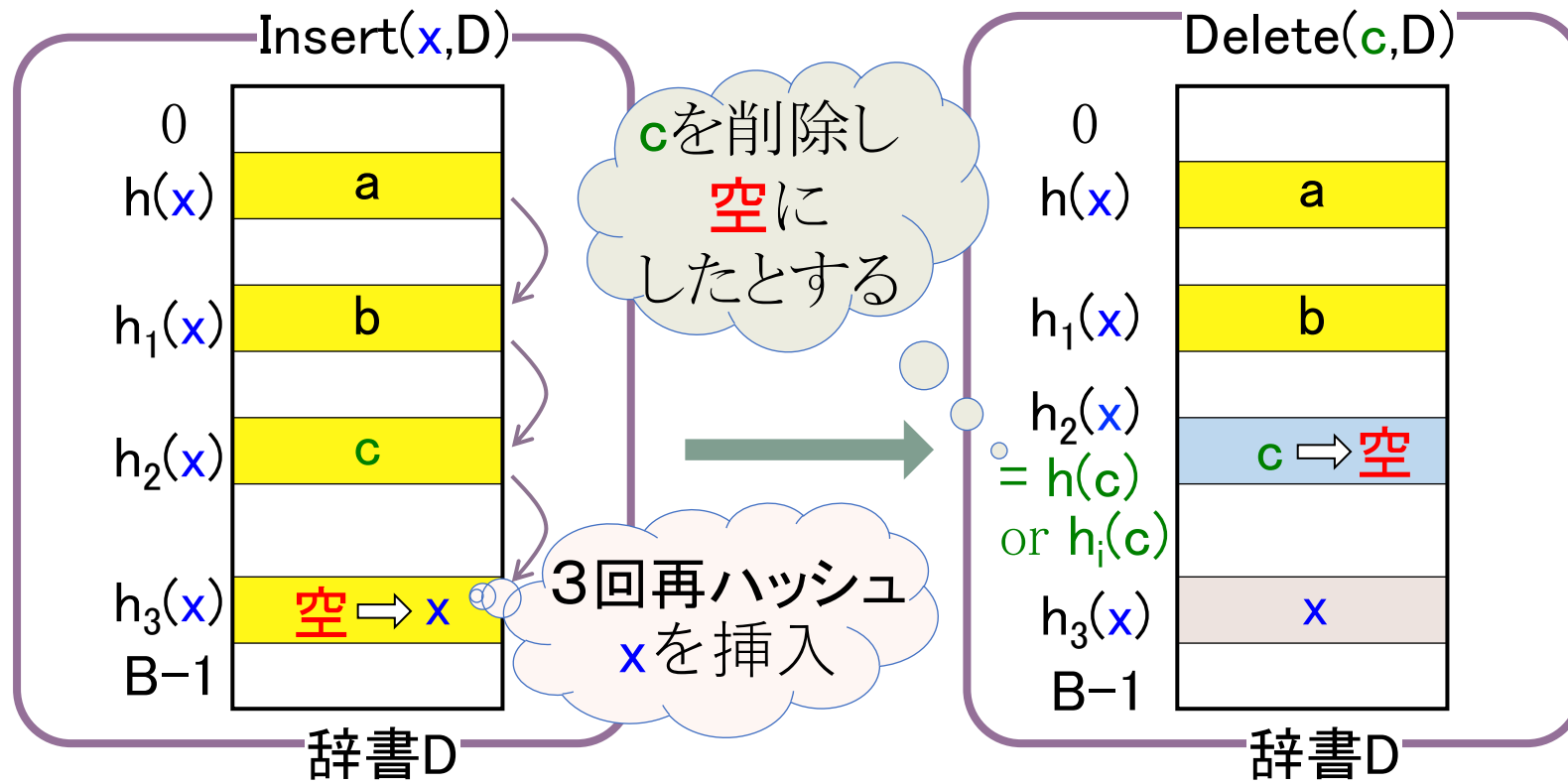
1. キーxをハッシュとするハッシュ値 $i = h(x)$ を求める
2. ハッシュ表の要素 $T[i]$ が
 - xならば、Member値を真 (1) にして、終了
 - 空の状態ならば、未登録なので、Member値を偽 (0) にし終了
 - x以外のキーならば、再ハッシュし再ハッシュ値 i を求め2.を繰り返す

キーxのD中での探索



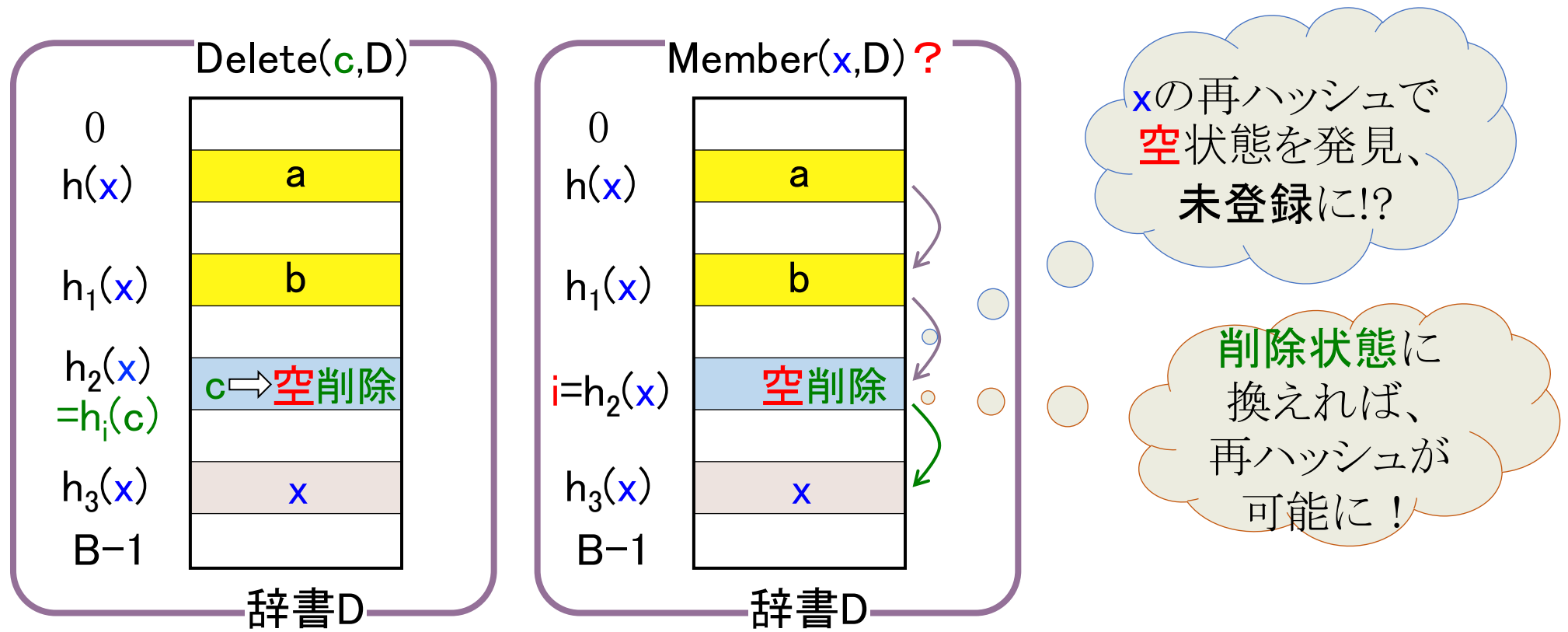
削除操作がある場合

- ハッシュ表の要素の状態
 - キー（辞書要素）でふさがっている**状態**
 - 一度でもキーでふさがったことのない**空状態**
 - 操作Deleteでキーが削除された**削除状態**



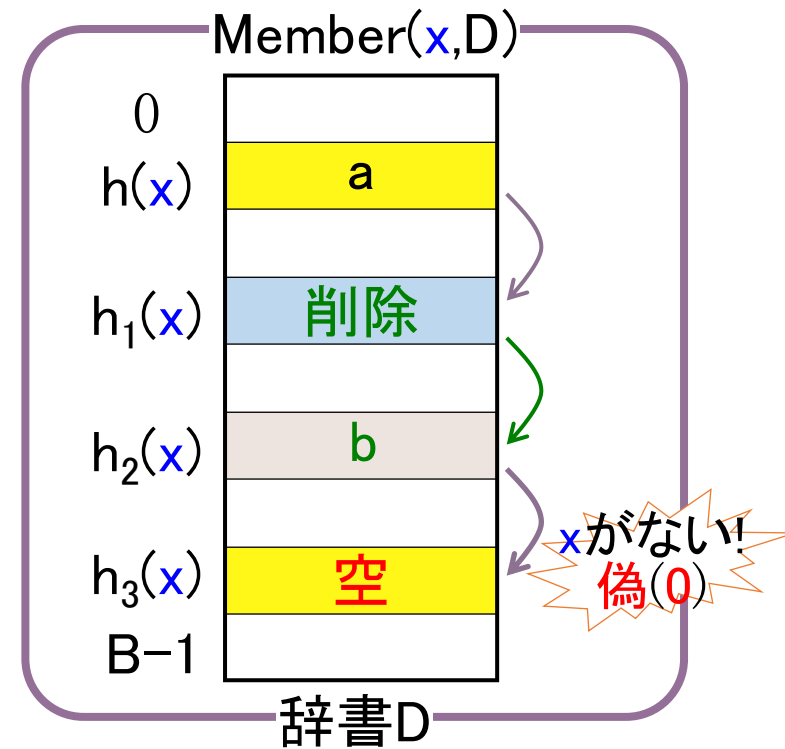
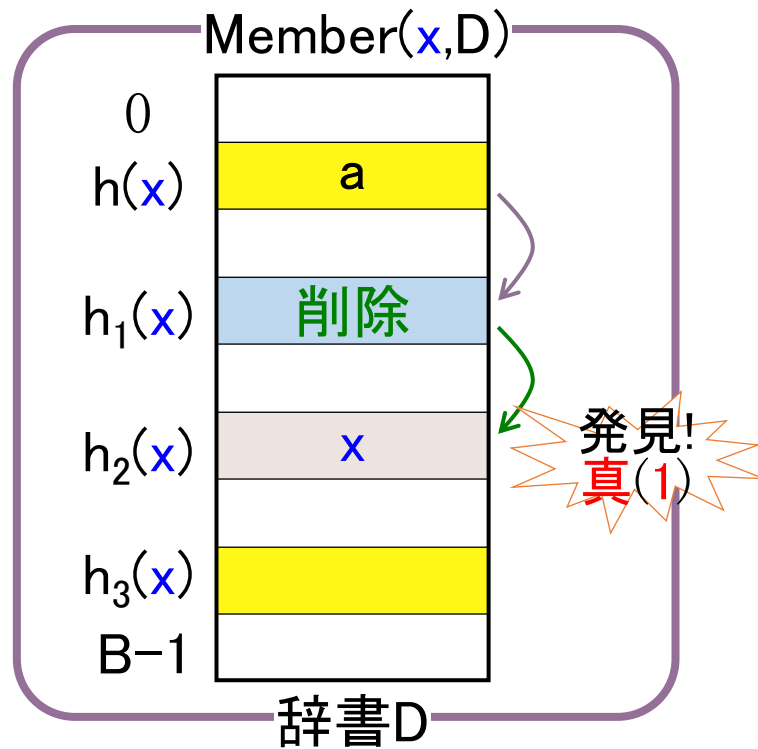
削除状態の必要性

- キー c を削除して、ハッシュ表の要素 $T[i]$ を **空状態** にすると
 - 他の操作（例えば、 x の Member 操作の再ハッシュ）で、 $T[i]$ にきたとき、**空状態** のため **キー x がない** と誤った判断をしてしまう
 - 削除後、**削除状態** にし、この先の **再ハッシュの必要性** を示す



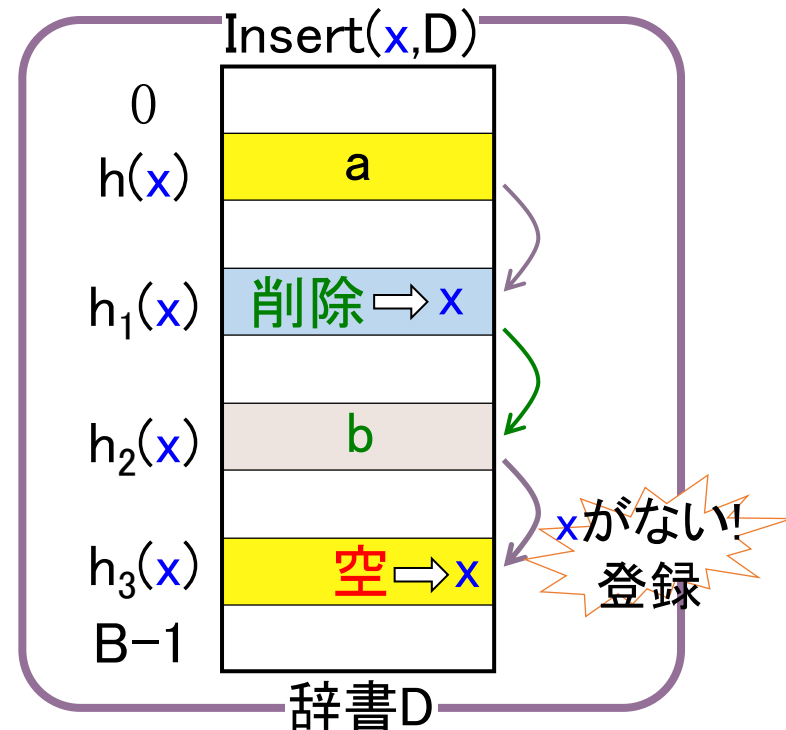
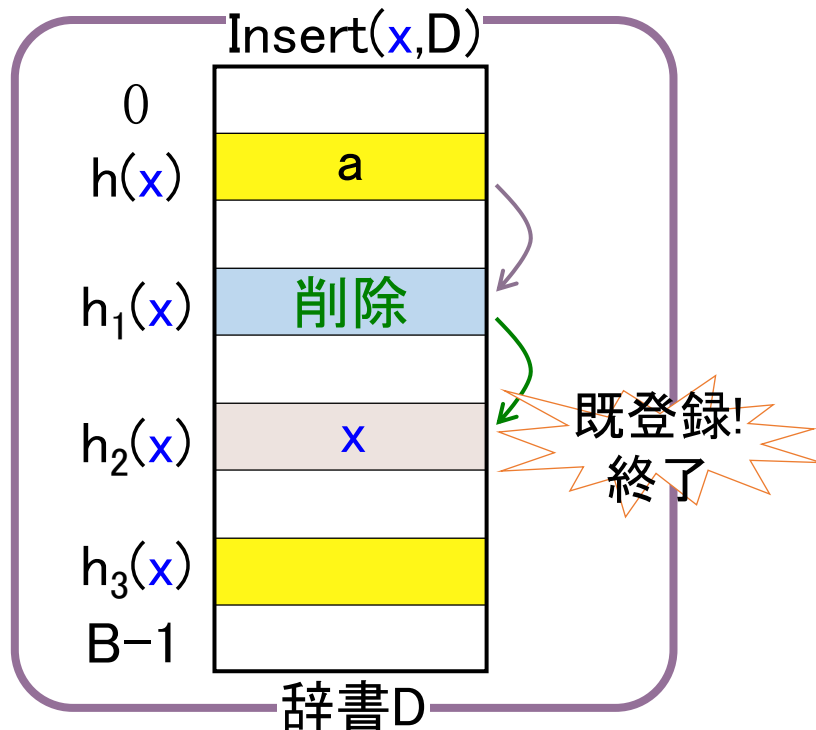
削除状態を設けた場合の所属Member(x,D)

1. キー x をハッシュとするハッシュ値 $i = h(x)$ を求める
2. ハッシュ表の要素 $T[i]$ が
 - x ならば, 既登録によりMember値を **真 (1)** にして, **終了**
 - **空** の状態ならば, **未登録** なので, Member値を **偽 (0)** にし **終了**
 - x 以外のキーまたは **削除状態** ならば, **再ハッシュ** し再ハッシュ値 i を求め2.を繰り返す



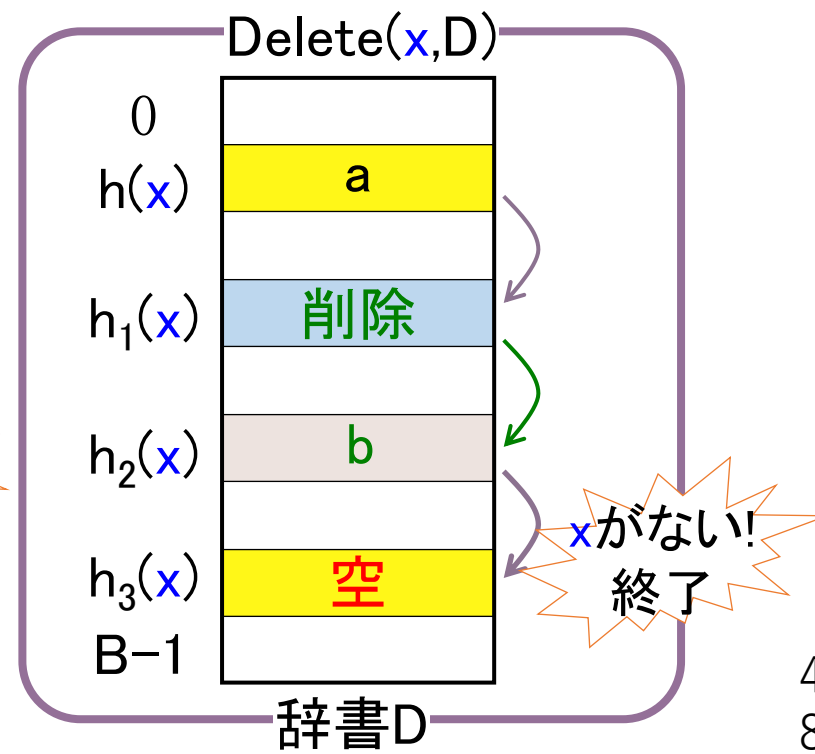
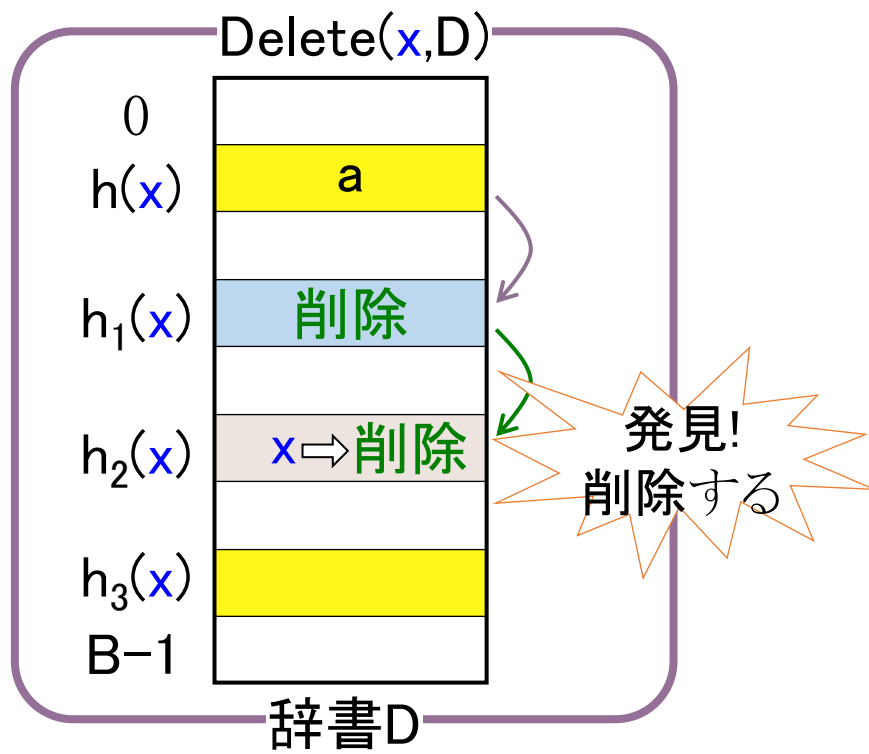
削除状態を設けた場合の挿入 $\text{Insert}(x, D)$

1. キー x をハッシュとするハッシュ値 $i = h(x)$ を求める
2. ハッシュ表の要素 $T[i]$ が
 - x ならば, **既登録** より **終了**
 - **空** の状態ならば, **未登録** なので, 3. にいき **登録**
 - x 以外のキーまたは **削除状態** ならば, **再ハッシュ** し再ハッシュ値 i を求め2. を繰り返す
3. **再ハッシュ** 中に **削除状態** の要素があれば, そこに x を **登録**. なければ, 最後の **空状態** に **登録**, **終了**



削除Delete(x, D)

1. キー x をハッシュとするハッシュ値 $i = h(x)$ を求める
2. ハッシュ表の要素 $T[i]$ が
 - x ならば, **削除状態**にして**終了**
 - **空**の状態ならば, **未登録**なので**終了**
 - x 以外 のキーまたは**削除状態**ならば, **再ハッシュ**し再ハッシュ値 i を求め2.を繰り返す



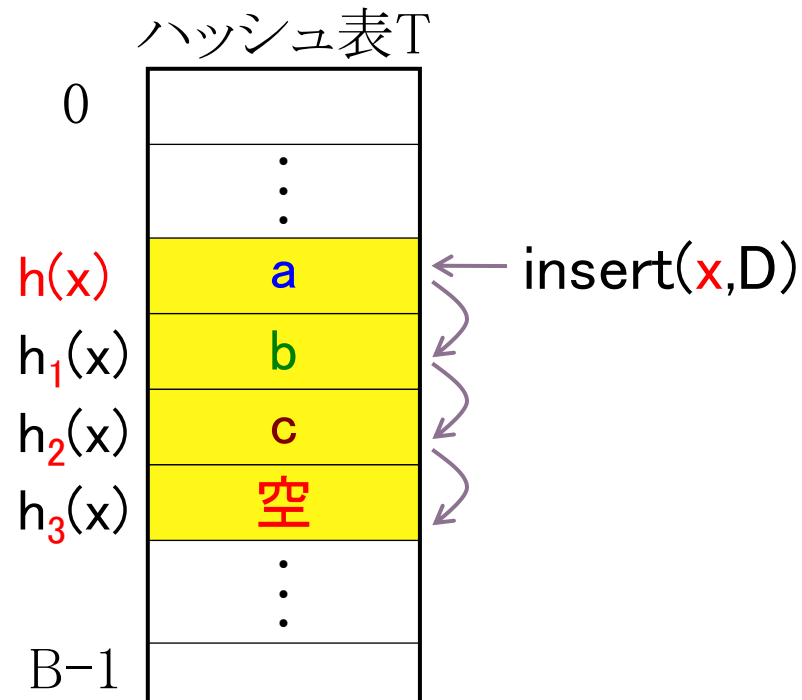
再ハッシュ関数の選定

- 線形検査法
- c 個離れた要素を調べる
- 2次関数検査法
- ランダムな順列を用いる

再ハッシュ関数：線形検査法

- $h_i(x) = (h(x) + i) \% B$ ($i = 1, \dots, B-1$), B はハッシュ表の大きさ
 - ハッシュ表を**環状**に考えて, **空**状態の要素が見つかるまで,
 $h(x) + 1, h(x) + 2, \dots$ と次の位置を調べていく

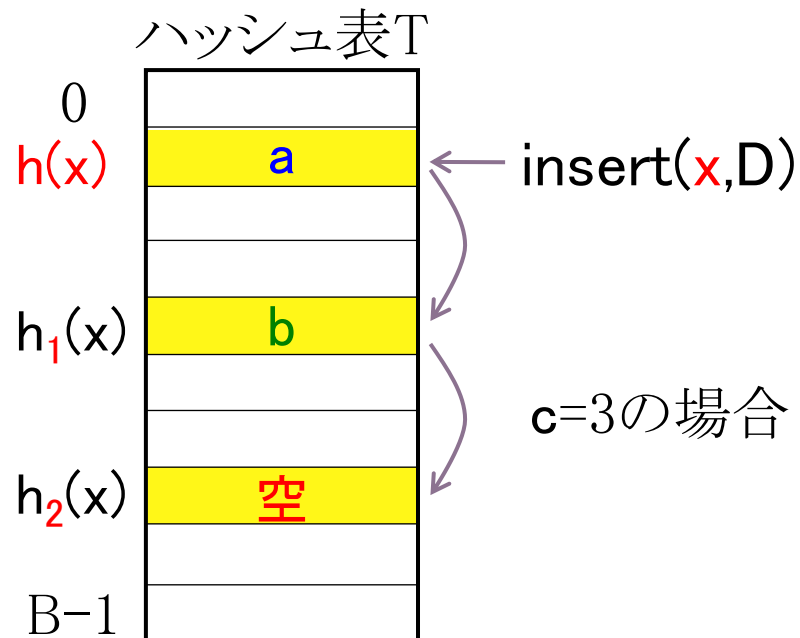
辞書要素が **$h(x)$** の**周り**に集中しやすい欠点を持つ



再ハッシュ関数：c個離れた要素

- $h_i(x) = (h(x) + ci) \% B$ ($i = 1, \dots, B-1$), c と B は**互いに素**
 - ハッシュ表を**環状**に考えて, **空**状態の要素が見つかるまで,
 $h(x) + c, h(x) + 2c, \dots$ と次の位置を調べていく

線形検索法と同じく, **c個目**ごとに**団子状態**になる欠点を持つ



再ハッシュ関数：2次関数検査法

- $h_i(x) = (h(x) + i^2) \% B \quad (i = 1, \dots, B-1)$
 - 団子状態は生じないが、ハッシュ表のすべてが検索されるとは限らない
 - Bが素数のときは、少なくともハッシュ表の半分は調べられる
 - 次の再帰的公式を用いれば、上述の再ハッシュ関数で、2乗演算を使わずにすむ
 - $h_{i+1}(x) = h_i(x) + d_i(x)$
 - $d_{i+1}(x) = d_i(x) + 2$
 - $h_0(x) = h(x)$
 - $d_0(x) = 1$

再ハッシュ関数：ランダムな順序

- $h_i(x) = (h(x) + d_i) \% B \quad (i = 1, \dots, B-1)$
 - d_1, d_2, \dots, d_{B-1} は、**ランダムな順列**で、シフトレジスタを用いて生成する方法がある

擬似乱数の生成方法

クローズドハッシュの平均時間計算量

- ハッシュ表の大きさ B
 - ハッシュ表は **N 個** の辞書要素でふさがっており, どのふさがり型も **等確率** とする
- ハッシュ表に **N 個登録** されている場合, **$N+1$ 個目** の新しい辞書要素を登録するときの平均時間計算量を考える
- **k 回目** に **空** 状態の要素を見つける **確率 P_k** を考える

$$P_1 = \frac{B-N}{B}$$

$$P_2 = \frac{N}{B} \frac{B-N}{B-1}$$

...

$$P_k = \frac{N}{B} \frac{N-1}{B-1} \cdots \frac{N-(k-2)}{B-(k-2)} \frac{B-N}{B-(k-1)}$$

N+1個目の辞書要素の登録

- **N+1個目**の新しい辞書要素をハッシュ表に**登録**するときの時間計算量を考える
 - **k回目**に**空**状態の要素を見つける確率 P_k

$$P_k = \frac{N}{B} \frac{N-1}{B-1} \cdots \frac{N-(k-2)}{B-(k-2)} \frac{B-N}{B-(k-1)}$$

〈N+1個目登録の平均時間計算量〉

$$= \sum_{k=1}^{N+1} k P_k = \frac{B+1}{B+1-N} \approx \frac{B}{B-N} = \frac{1}{1-\alpha}$$

ただし、 $\alpha = N/B$ で、
ハッシュ表中の辞書要素の割合

1個登録に要する平均時間計算量

- ハッシュ表にN個まで登録したとき、1個登録に要した平均時間計算量は、

〈N個まで登録したときの1個登録平均時間計算量〉

$$= \frac{1}{N} \sum_{k=1}^N \langle k \text{ 個目登録の平均時間計算量} \rangle$$

$$= \frac{1}{N} \sum_{k=1}^N \frac{B+1}{B+2-k} = \frac{B+1}{N} \ln \frac{B+1}{B+2-N} \doteq \frac{B}{N} \ln \frac{B}{B-N}$$

$$= \frac{1}{\alpha} \ln \frac{1}{1-\alpha} \quad \text{ただし、}\alpha = N/B \text{で、}$$

ハッシュ表中の辞書要素の割合

- ◆ ハッシュ表の90% ($\alpha = 0.9$) を辞書要素でうめたときの1個登録に要した平均時間計算量は、2.56。

探索に要する平均時間計算量

- ハッシュ表に**N個まで**登録したとき

〈表にない要素を探す平均時間計算量〉

= 〈N + 1 個目登録の平均時間計算量〉

$$= \frac{1}{1 - \alpha}$$

ただし、 $\alpha = N/B$ で、
ハッシュ表中の辞書要素の割合

〈表にある要素を探す平均時間計算量〉

= 〈N個まで登録したときの1個登録平均時間計算量〉

$$= \frac{1}{\alpha} \ln \frac{1}{1 - \alpha}$$

削除に要する平均時間計算量

- ハッシュ表にN個まで登録したとき

〈要素を削除する平均時間計算量〉

= 〈その要素を探索する平均時間計算量〉

クローズドハッシュの平均時間計算量まとめ

- 登録要素数 N がハッシュ表の90%以上（ $\alpha > 0.9$ ）のとき
 - 表中にある要素を探す平均時間計算量は**2.56以上**
 - 表中にない要素を探す平均時間計算量は**10以上**
- 表の再構成
 - 登録要素数がハッシュ表の**90%以上**を占めたとき
 - 大きさが2倍の**新しいハッシュ表**を作る
 - ⇒それぞれの平均時間計算量が**2.56以下**および**10以下**に保つことはできる
 - ⇒登録要素数 **N** は表の大きさ **B** の**90%**程度
- 最小の要素を取り出す操作**Min**の効率のよいアルゴリズムはない

まとめ

- 集合
 - 集合の仕様
 - 実現：ビットベクトル
 - 実現：連結リスト
- 辞書
 - 辞書の仕様
 - 実現：配列へのベタ詰め法
 - 実現：ハッシュ法
 - オープンハッシュ法
 - クローズドハッシュ法

演習1 オープンハッシュ法

- 抽象データ型辞書をオープンハッシュ法で実現する.
 - 辞書操作, 挿入, 探索(所属), 削除について以下の設問に答えよ.
 - ただし,
 - ハッシュ表の大きさBは5, インデックスは0からB-1までとする.
 - 辞書要素(キー)の集合 = $\{0, 1, \dots, 14\}$ とする.
 - ハッシュ関数 $h(x) = x \% B$ とする.
 - 最初、ハッシュ表(辞書)には何も登録されていないとする.
1. 5つの辞書要素(キー) 2, 8, 14, 3, 9 をこの順でハッシュ表に登録する. このときのハッシュ表への登録手順を概説し, できあがったハッシュ表の図を描け.
 2. つづいて、9, 13がハッシュ表に登録されているか否かを調べる. このときのハッシュ表の探索手順を概説せよ.
 3. 最後に、9, 13をハッシュ表から削除する. このときのハッシュ表の削除手順を概説し, できあがったハッシュ表の図を描け.

演習2 クローズドハッシュ法

- 抽象データ型辞書Dをクローズドハッシュ法で実現する.
- ただし,
 - 辞書Dのハッシュ表Tの大きさBは11, そのインデックスは0からB-1とする.
 - 辞書要素(キー)の集合 = $\{0, 1, \dots, 109\}$ とする.
 - ‘空の状態’を -1, ‘削除状態’を -2 で表す.
 - ハッシュ関数 $h(x) = x \% B$ とする.
 - 再ハッシュ関数 $h_i(x) = (h(x) + 3 * i) \% B$ とする.
 - いま、辞書D(ハッシュ表T)には右図のように、キーが登録されており、以下の設問を順に操作していくとする.

| ハッシュ表T | |
|--------|----|
| 0 | -1 |
| 1 | -1 |
| 2 | 32 |
| 3 | -1 |
| 4 | -1 |
| 5 | -1 |
| 6 | -1 |
| 7 | 18 |
| 8 | -1 |
| 9 | -1 |
| 10 | 21 |

1. 辞書要素(キー) 7を登録する。そのときのハッシュ, 再ハッシュの手順を説明し, 結果のハッシュ表Tを図示せよ.
2. 辞書要素(キー) 32を削除する。そのときのハッシュ, 再ハッシュの手順を説明し, 結果のハッシュ表Tを図示せよ.
3. 再び、辞書要素(キー) 7を登録するときのハッシュ, 再ハッシュの手順を説明し, 結果のハッシュ表Tを図示せよ.

提出方法

- ドローイングソフトを使ってもかまいませんが、手書きを写真でとったものでOKです.
- pdfや画像フォーマットで提出してください
- 提出方法：LETUS
- 締め切り：2023/7/10 10:30まで