

# CHAPTER 8

## HASHING

All the programs in this file are selected from

Ellis Horowitz, Sartaj Sahni, and Susan Anderson-Freed  
“Fundamentals of Data Structures in C”,

# Symbol Table

- Definition

  - A set of name-attribute pairs

- Operations

  - Determine if a particular name is in the table
  - Retrieve the attributes of the name
  - Modify the attributes of that name
  - Insert a new name and its attributes
  - Delete a name and its attributes

# The ADT of Symbol Table

Structure SymbolTable(SymTab) is

objects: a set of name-attribute pairs, where the names are unique

functions:

for all *name* belongs to *Name*, *attr* belongs to *Attribute*, *symtab* belongs to *SymbolTable*, *max\_size* belongs to integer

SymTab Create(*max\_size*) ::= create the empty symbol table whose maximum capacity is *max\_size*

Boolean IsIn(*symtab*, *name*) ::= if (*name* is in *symtab*) return TRUE  
else return FALSE

Attribute Find(*symtab*, *name*) ::= if (*name* is in *symtab*) return the corresponding attribute else return null attribute

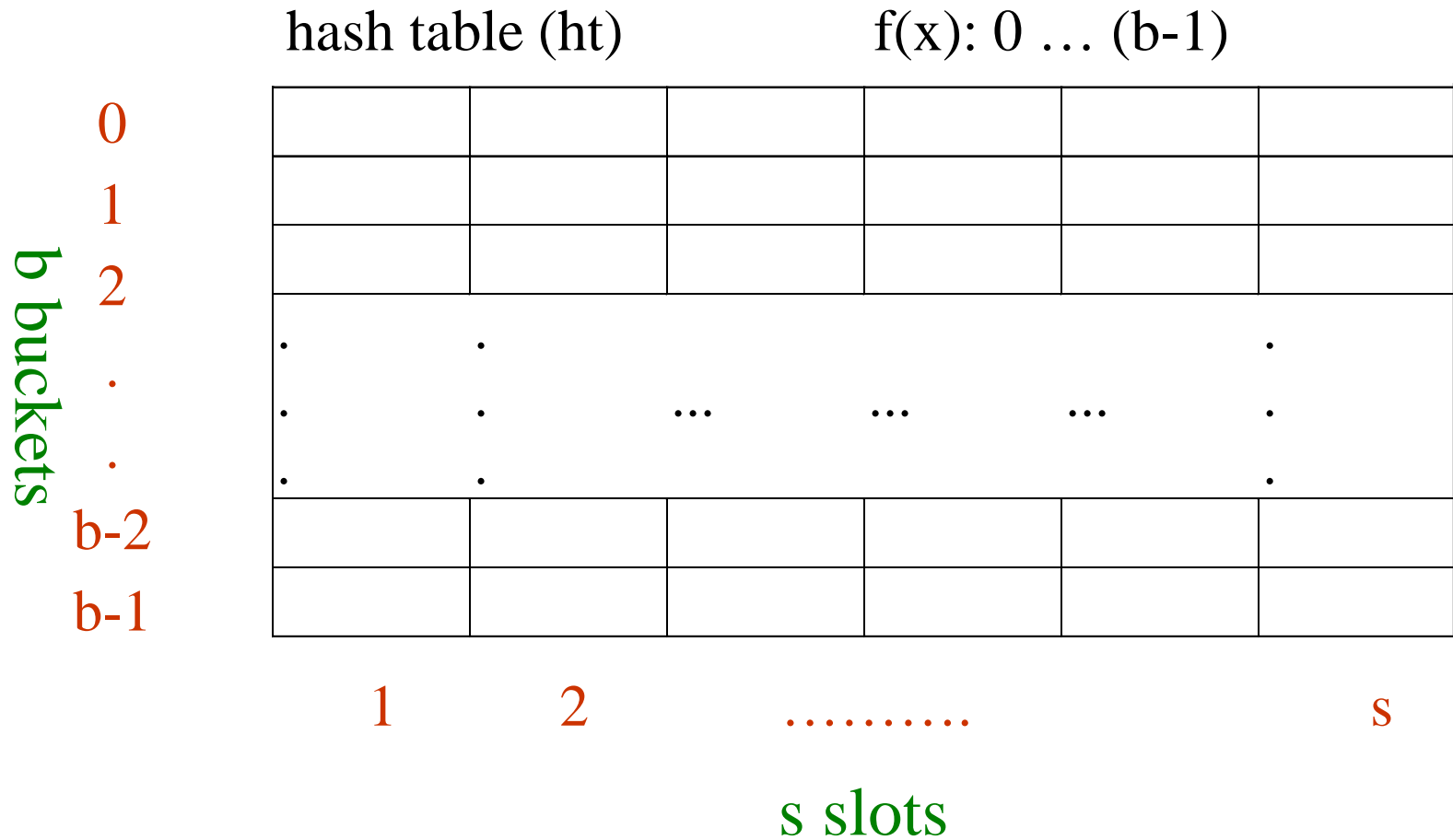
SymTab Insert(*symtab*, *name*, *attr*) ::= if (*name* is in *symtab*)  
replace its existing attribute with *attr*  
else insert the pair (*name*, *attr*) into  
*symtab*

SymTab Delete(*symtab*, *name*) ::= if (*name* is not in *symtab*) return  
else delete (*name*, *attr*) from *symtab*

# Search vs. Hashing

- Search tree methods: **key comparisons**
- Hashing methods: **hash functions**
- types
  - statistic hashing
  - dynamic hashing

# Static Hashing





# Identifier Density and Loading Density

- The *identifier density* of a hash table is the ratio  $n/T$ 
  - $n$  is the number of identifiers in the table
  - $T$  is possible identifiers
- The *loading density* or *loading factor* of a hash table is  $\alpha = n/(sb)$ 
  - $s$  is the number of slots
  - $b$  is the number of buckets



# Synonyms

- Since the number of buckets  $b$  is usually several orders of magnitude lower than  $T$ , the hash function  $f$  must map several different identifiers into the same bucket
- Two identifiers,  $i$  and  $j$  are **synonyms** with respect to  $f$  if  $f(i) = f(j)$

# Overflow and Collision

- An **overflow** occurs when we hash a new identifier into a full bucket
- A **collision** occurs when we hash two non-identical identifiers into the same bucket





# Example

	Slot 0	Slot 1
0	acos	atan <span style="color: red;">synonyms</span>
1		
2	char	ceil
3	define	
4	exp	
5	float	floor <span style="color: red;">synonyms</span>
6		
...		
25		

synonyms:  
char, ceil,  
clock, ctime



overflow

$b=26, s=2, n=10, \alpha=10/52=0.19, f(x)=\text{the first char of } x$

$x$ : acos, define, float, exp, char, atan, ceil, floor, clock, ctime

$f(x)$ : 0, 3, 5, 4, 2, 0, 2, 5, 2, 2

# Hashing Functions

- Two requirements
  1. easy computation
  2. minimal number of collisions
- mid-square (middle of square)

$$f_m(x) = \text{middle}(x^2)$$

- division

$$f_D(x) = x \% M \quad (0 \sim (M-1))$$

Avoid the choice of M that leads to many collisions



# Hashing Functions

## ■ Folding

- Partition the identifier  $x$  into several parts
- All parts except for the last one have the same length
- Add the parts together to obtain the hash address
- $K=12320324111220$
- Two possibilities
  - **Shift folding**
    - $x_1=123, x_2=203, x_3=241, x_4=112, x_5=20, \text{address}=699$
  - **Folding at the boundaries**
    - $x_1=123, x_2=203, x_3=241, x_4=112, x_5=20, \text{address}=897$

shift folding

123	203	241	112	20
P1	P2	P3	P4	P5

123 →

203 →

241 →

112 →

20

699

folding at  
the boundaries

123

203

241

112

20

→

←

→

←

→

MSD ---> LSD

LSD <--- MSD

302

211

897

# Digital Analysis

- All the identifiers are known in advance

$M=1\sim 999$

$X_1$	$d_{11}$	$d_{12}$	$\dots$	$d_{1n}$
$X_2$	$d_{21}$	$d_{22}$	$\dots$	$d_{2n}$
$\dots$				
$X_m$	$d_{m1}$	$d_{m2}$	$\dots$	$d_{mn}$

Select 3 digits from  $n$

Criterion:

Delete the digits having the most skewed distributions

# Overflow Handling

- Linear Open Addressing (linear probing)
- Quadratic probing
- Chaining

# Data Structure for Hash Table

```
#define MAX_CHAR 10
#define TABLE_SIZE 13
typedef struct {
    char key[MAX_CHAR];
    /* other fields */
} element;
element hash_table[TABLE_SIZE];
```

# Example

Identifier	Additive Transform	x	Hash
for	102+111+114	327	2
do	100+111	211	3
while	119+104+105+108	537	4
if	105+102	207	12
else	101+108+115+101	425	9
function	102+117+110+99+	870	12

[0] ↵ **function** ↵ ↵

[1] ↵ ↵ ↵

[2] ↵ **for** ↵ ↵

[3] ↵ **do** ↵ ↵

[4] ↵ **while** ↵ ↵

[5] ↵ ↵ ↵

[6] ↵ ↵ ↵

[7] ↵ ↵ ↵

[8] ↵ ↵ ↵

[9] ↵ **else** ↵ ↵

[10] ↵ ↵ ↵

[11] ↵ ↵ ↵

[12] ↵ **if** ↵ ↵

[12] 英



# Linear Probing (linear open addressing)

- Compute  $f(x)$  for identifier  $x$

- Examine the buckets

$ht[(f(x)+j)\%TABLE\_SIZE]$

$0 \leq j \leq TABLE\_SIZE$

- The bucket contains  $x$ .

- The bucket contains the empty string

- Return to  $ht[f(x)]$

桶	x	桶的搜尋次數
0	acos	1
1	atoi	2
2	char	1
3	define	1
4	exp	1
5	ceil	4
6	cos	5
7	float	3
8	atol	9
9	floor	5
10	ctime	9
...		
25		

# Linear Probing

```
void linear_insert(element item, element ht[])
{
    int i, hash_value;
    i = hash_value = hash(item.key);
    while(strlen(ht[i].key)) {
        if (!strcmp(ht[i].key, item.key))
            fprintf(stderr, "Duplicate entry\n");
            exit(1);
        }
        i = (i+1)%TABLE_SIZE;
        if (i == hash_value) {
            fprintf(stderr, "The table is full\n");
            exit(1);
        }
    }
    ht[i] = item;
}
```



# Problem of Linear Probing

- Identifiers tend to cluster together
- Adjacent cluster tend to coalesce
- Increase the search time



# Quadratic Probing

- Linear probing searches buckets  $(f(x)+i)\%b$
- Quadratic probing uses a quadratic function of  $i$  as the increment
- Examine buckets  $f(x)$ ,  $(f(x)+i^2)\%b$ ,  $(f(x)-i^2)\%b$ , for  $1\leq i\leq(b-1)/2$
- $b$  is a prime number of the form  $4j+3$ ,  $j$  is an integer

# Rehashing

- Try  $f_1, f_2, \dots, f_m$  in sequence if collision occurs
- disadvantage
  - comparison of identifiers with different hash values
  - use chain to resolve collisions

# Data Structure for Chaining

```
#define MAX_CHAR 10
#define TABLE_SIZE 13
#define IS_FULL(ptr) (!ptr)
typedef struct {
    char key[MAX_CHAR];
    /* other fields */
} element;
typedef struct list *list_pointer;
typedef struct list {
    element item;
    list_pointer link;
};
list_pointer hash_table[TABLE_SIZE];
```

# Chain Insert

```
void chain_insert(element item, list_pointer ht[])
{
    int hash_value = hash(item.key);
    list_pointer ptr, trail=NULL, lead=ht[hash_value];
    for (; lead; trail=lead, lead=lead->link)
        if (!strcmp(lead->item.key, item.key)) {
            fprintf(stderr, "The key is in the table\n");
            exit(1);
        }
    ptr = (list_pointer) malloc(sizeof(list));
    if (IS_FULL(ptr)) {
        fprintf(stderr, "The memory is full\n");
        exit(1);
    }
    ptr->item = item;
    ptr->link = NULL;
    if (trail) trail->link = ptr;
    else ht[hash_value] = ptr;
}
```



# Results of Hash Chaining

acos, atoi, char, define, exp, ceil, cos, float, atol, floor, ctime  
f(x)=first character of x

[0]	-> acos -> atoi -> atol
[1]	-> NULL
[2]	-> char -> ceil -> cos -> ctime
[3]	-> define
[4]	-> exp
[5]	-> float -> floor
[6]	-> NULL
...	
[25]	-> NULL

# of key comparisons= $21/11=1.91$   
CHAPTER 8

The *loading density* or *loading factor* of a hash table is  
 $\alpha = n/(sb)$

$\alpha=n/b$	.50	.75	.90	.95
hashing function	chain/open	chain/open	chain/open	chain/open
mid square	1.26/1.73	1.40/9.75	1.45/37.14	1.47/37.53
division	1.19/4.52	1.31/7.20	1.38/22.42	1.41/25.79
shift fold	1.33/21.75	1.48/65.10	1.40/77.01	1.51/118.57
Bound fold	1.39/22.97	1.57/48.70	1.55/69.63	1.51/97.56
digit analysis	1.35/4.55	1.49/30.62	1.52/89.20	1.52/125.59
theoretical	1.25/1.50	1.37/2.50	1.45/5.50	1.48

$n$  is the number of identifiers in the table

$b$  is the number of buckets

# Dynamic Hashing (extensible hashing)

- Dynamically increasing and decreasing file size
- concepts
  - file: a collection of records
  - record: a key + data, stored in pages (buckets)
  - space utilization

$$\frac{\text{Number of Record}}{\text{Number of Pages} * \text{Page Capacity}}$$

# Dynamic Hashing Using Directories

Example.  $m(\# \text{ of pages})=4$ ,  $P(\text{page capacity})=2$

00, 01, 10, 11

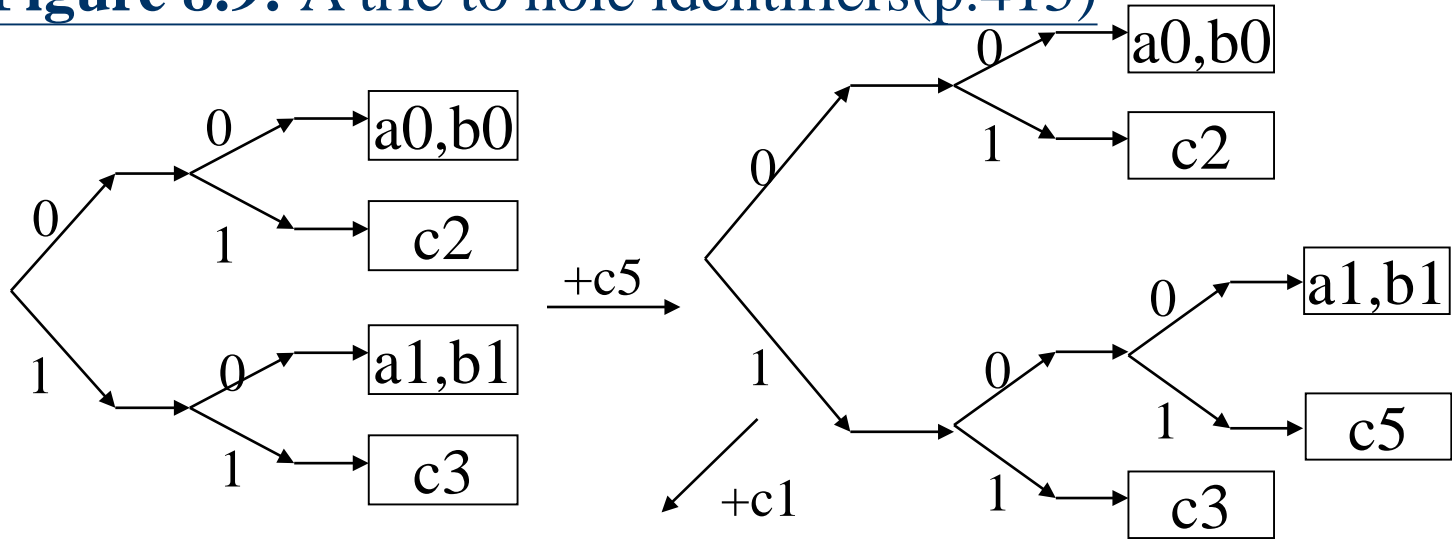
Identifiers	Binary representaiton
a0	100 <u>000</u>
a1	100 <u>001</u>
b0	101 <u>000</u>
b1	101 <u>001</u>
c0	110 <u>000</u>
c1	110 <u>001</u>
c2	110 <u>010</u>
c3	110 <u>011</u>

from LSB  
to MSB

allocation:  
lower order  
two bits

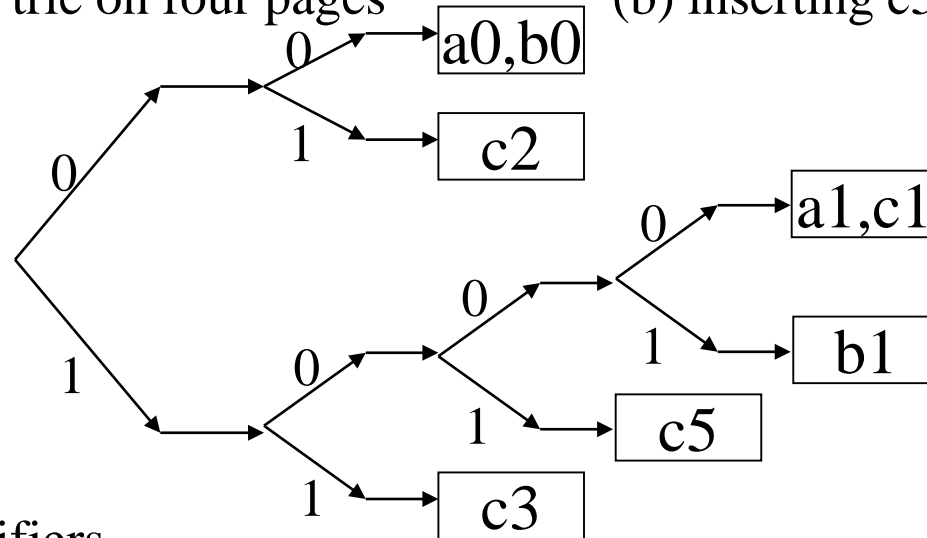
**\*Figure 8.8:** Some identifiers requiring 3 bits per character(p.414)

**\*Figure 8.9: A trie to hole identifiers(p.415)**



(a) two level trie on four pages

(b) inserting c5 with overflow



(c) inserting c1 with overflow

Note: time to access  
a page: # of bits to  
distinguish the identifiers  
Note: identifiers skewed:  
depth of tree skewed

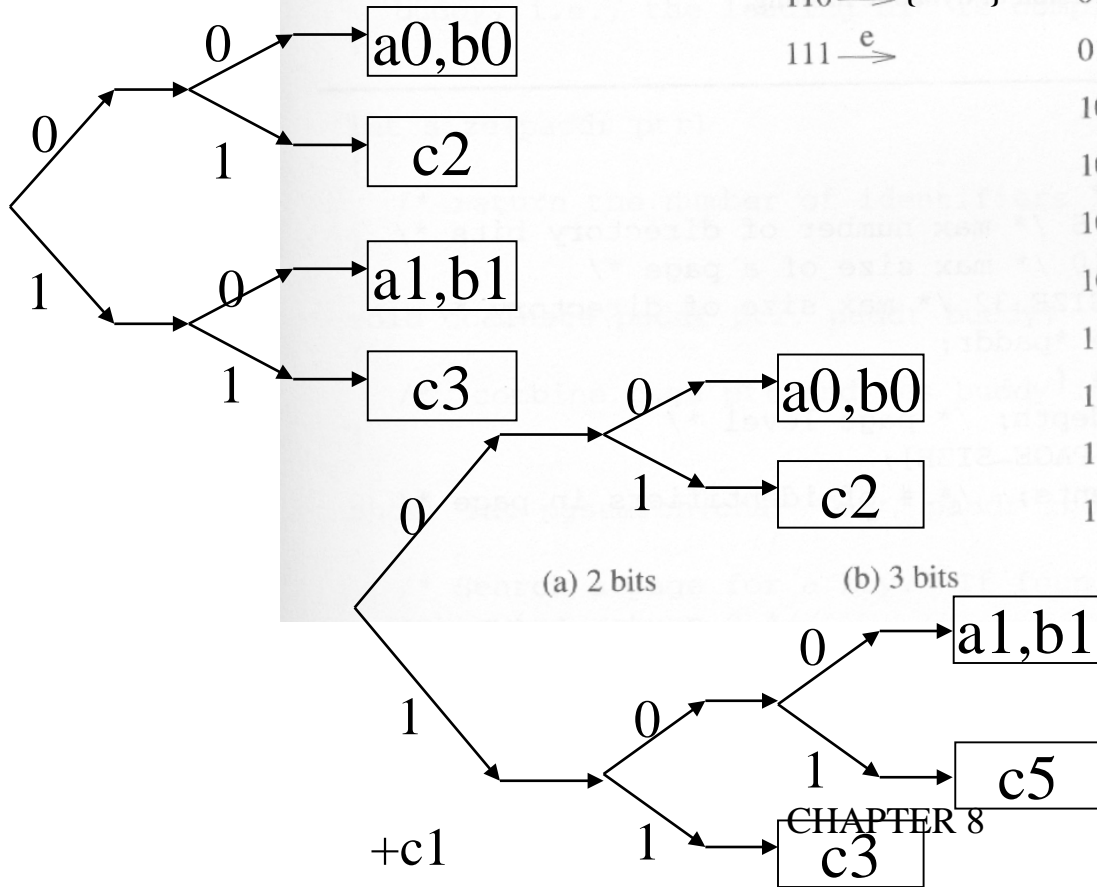
# Extendible Hashing

local depth  
global depth: 4

$f(x)$  = a set of binary digits --> table lookup

00 → a	a0, b0	000 → a	a0, b0	0000 → a	2	{0000,1000,0100,1100}
01 → c	a1, b1	001 → c	{000,100}	0001 → c	4	{0001}
10 → b	c2	010 → b	{001}	0010 → b	2	{0010,1010,0110,1110}
11 → d	c3	011 → e	{010,110}	0011 → f	2	{0011,1011,0111,1111}
		100 → a	{011,111}	0100 → a		
		101 → d	c5	0101 → e	3	{0101,1101}
		110 → b	{101}	0110 → b		
		111 → e		0111 → f		
				1000 → a		
				1001 → d	4	b1 {1001}
				1010 → b		
				1011 → f		
				1100 → a		
				1101 → e		
				1110 → b		
				1111 → f		

page pointer



pages c & d: buddies

If keys do not uniformly divide up among pages, then the directory can grow quite large, but most of entries will point to the same page

f: a family of hashing functions

$\text{hash}_i$ : key  $\rightarrow \{0 \dots 2^{i-1}\}$   $1 \leq i \leq d$

$\text{hash}(\text{key}, i)$ : produce random number of  $i$  bits from identifier key

$\text{hash}_i$  is  $\text{hash}_{i-1}$  with either a zero or one appeared as the new leading bit of result

<u>100 000</u> $\text{hash}(a0,2)=00$	<u>100 001</u> $\text{hash}(a1,4)=0001$	<u>101 000</u> $\text{hash}(b0,2)=00$
<u>101 001</u> $\text{hash}(b1,4)=1001$	<u>110 001</u> $\text{hash}(c1, 4)=0001$	<u>110 010</u> $\text{hash}(c2, 2)=10$
<u>110 011</u> $\text{hash}(c3,2)=11$	<u>110 101</u> $\text{hash}(c5,3)=101$	

## \*Program 8.5: Dynamic hashing (p.421)

```
#include <stdio.h>
#include <alloc.h>
#include <stdlib.h>    25=32
#define WORD_SIZE 5    /* max number of directory bits */
#define PAGE_SIZE 10  /* max size of a page */
#define DIRECTORY_SIZE 32 /* max size of directory */
typedef struct page *paddr;
typedef struct page {
    int local_depth; /* page level */
    char *name[PAGE_SIZE]; /* the actual identifiers */
    int num_idents; /* #of identifiers in page */
};
typedef struct {
    char *key; /* pointer to string */
    /*other fields */
} brecord;
int global_depth; /* trie height */
paddr directory[DIRECTORY_SIZE]; /* pointers to pages */
```

See Figure 8.10(c) global depth=4  
local depth of a =2



```
paddr hash(char *, short int);
paddr buddy(paddr);
short int pgsearch(char *, paddr );
int convert(paddr);
void enter(brecord, paddr);
void pgdelete(char *, paddr);
paddr find(brecord, char *);
void insert (brecord, char *);
int size(paddr);
void coalesce (paddr, paddr);
void delete(brecord, char *);
```

```
paddr hash(char *key, short int precision)
```

```
{
```

```
    /* *key is hashed using a uniform hash function, and the
       low precision bits are returned as the page address */
```

```
}
```

directory subscript for directory lookup

```

paddr buddy(paddr index)
{
    /*Take an address of a page and returns the page's
    buddy, i. e., the leading bit is complemented */
    buddy  $b_{n-1}b_{n-2} \dots b_0$        $b_{n-1}b_{n-2} \dots b_0$ 
}

int size(paddr ptr)
{
    /* return the number of identifiers in the page */
}

void coalesce(paddr ptr, paddr, buddy)
{
    /*combine page ptr and its buddy into a single page */
}

short int pgsearch(char *key, paddr index)
{
    /*Search a page for a key. If found return 1
    otherwise return 0 */
}

```

```
void convert (paddr ptr)
{
    /* Convert a pointer to a pointer to a page to an equivalent integer */
}
```

```
void enter(brecord r, paddr ptr)
{
    /* Insert a new record into the page pointed at by ptr */
}
```

```
void pgdelete(char *key, paddr ptr)
{
    /* remove the record with key, hey, from the page pointed to by ptr */
}
```

```
short int find (char *key, paddr *ptr)
{
    /* return 0 if key is not found and 1 if it is. Also,
       return a pointer (in ptr) to the page that was searched.
       Assume that an empty directory has one page. */
}
```

```

paddr index;
int intindex;
index = hash(key, global_depth);
intindex = convert(index);
*ptr = directory[intindex];
return pgsearch(key, ptr);
}

```

```

void insert(brecord r, char *key)

```

```

{
    paddr ptr;
    if find(key, &ptr) {
        fprintf(stderr, " The key is already in the table.\n");
        exit(1);
    }
    if (ptr-> num_idents != PAGE_SIZE) {
        enter(r, ptr);
        ptr->num_idents++;
    }
}

```

```

else{ /*Split the page into two, insert the new key, and update global_depth
        if necessary.

```

```
If this causes global_depth to exceed WORD_SIZE then print an error
and terminate. */
```

```
};
}
```

```
void delete(brecord r, char *key)
```

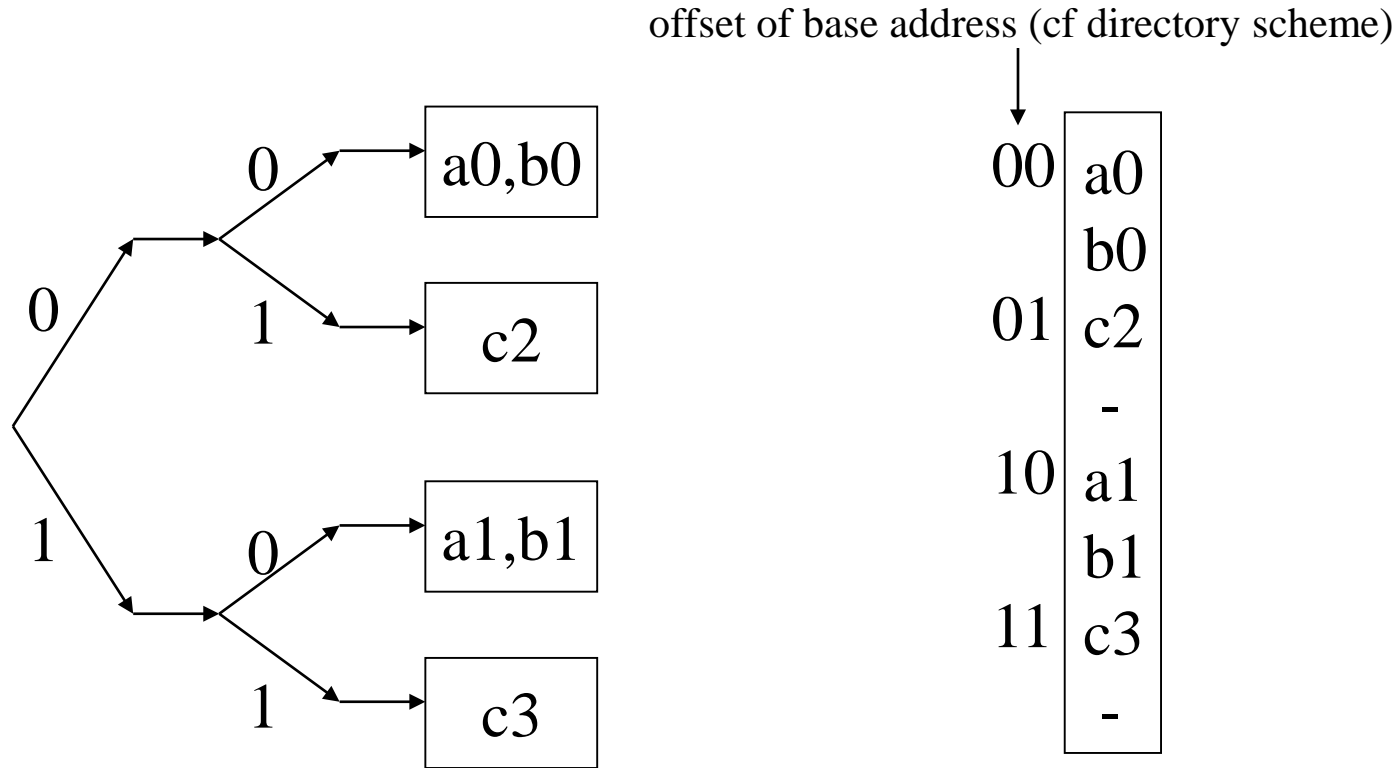
```
{
/* find and delete the record r from the file */
  paddr ptr;
  if (!find (key, &ptr )) {
    fprintf(stderr, "Key is not in the table.\n");
    return; /* non-fatal error */
  }
  pgdelete(key, ptr);
  if (size(ptr) + size(buddy(ptr)) <= PAGE_SIZE)
    coalesce(ptr, buddy(ptr));
}
```

```
void main(void)
```

```
{
}
```

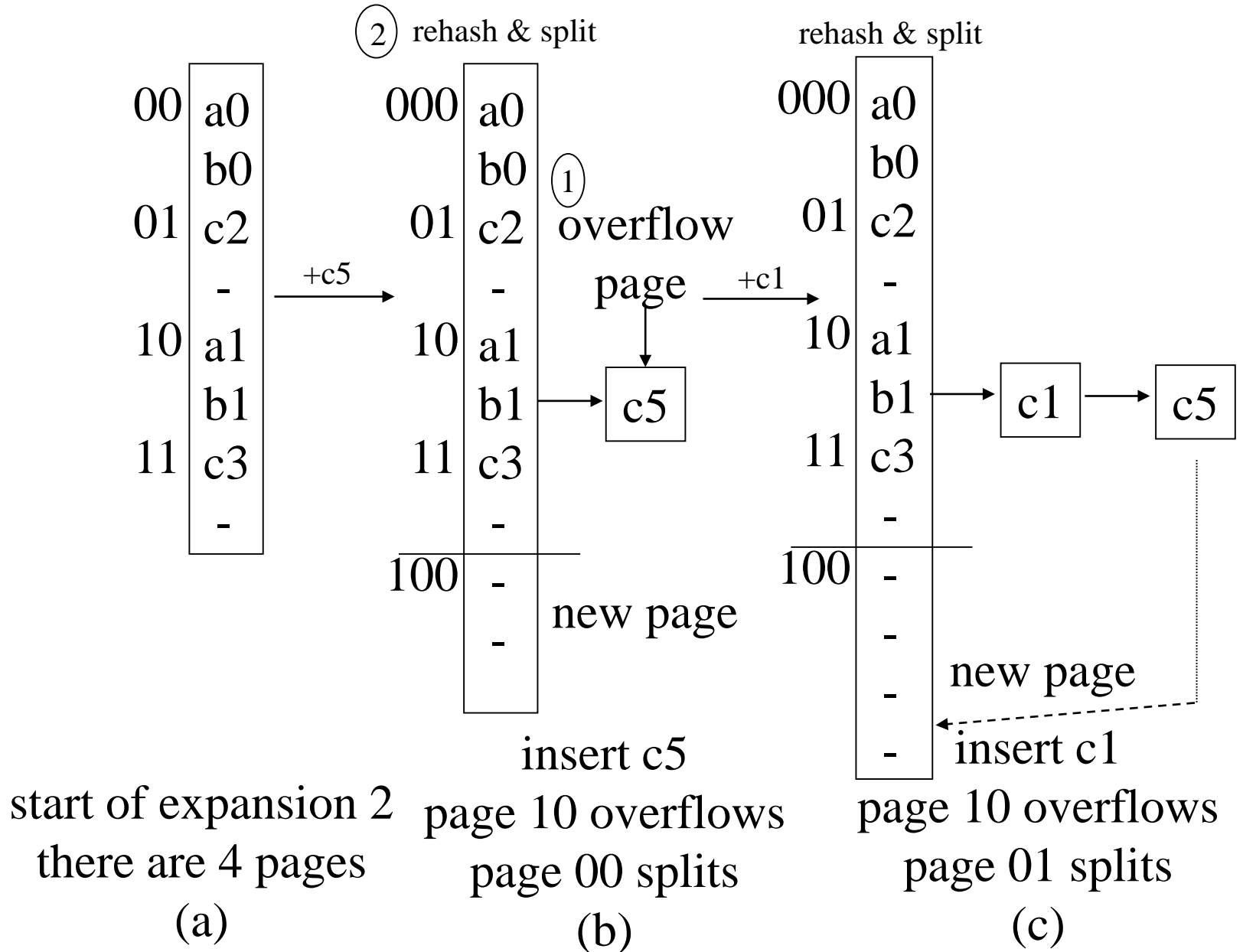
# Directoryless Dynamic Hashing (Linear Hashing)

continuous address space

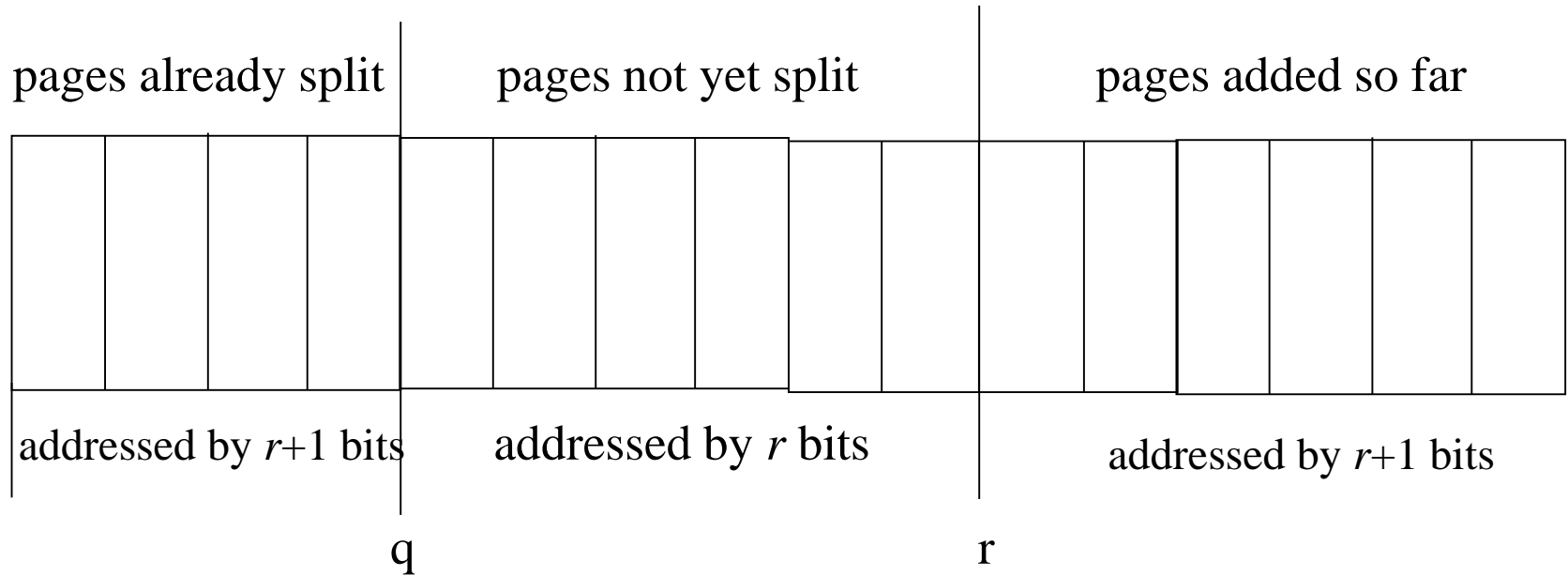


\*Figure 8.12: A trie mapped to a directoryless, contiguous storage (p.424)

**\*Figure 8.13: An example with two insertions (p.425)**



**\*Figure 8.14:** During the  $r$ th phase of expansion of directoryless method (p.426)



←  $2^r$  pages at start →

suppose we are at phase  $r$ ; there are  $2^r$  pages indexed by  $r$  bits



**\*Program 8.6: Modified hash function (p.427)**

```
if ( hash(key,r) < q)
    page = hash(key, r+1);
else
    page = hash(key, r);
if needed, then follow overflow pointers;
```