



# Hash Tables with External Chaining

by Andrew W. Appel and Robert M. Dondero Jr.,  
Princeton University



# Key-value store

## Maintain a collection of key/value pairs

- Each key is a **string**; each value is an **int**
- Unknown number of key-value pairs

## Examples

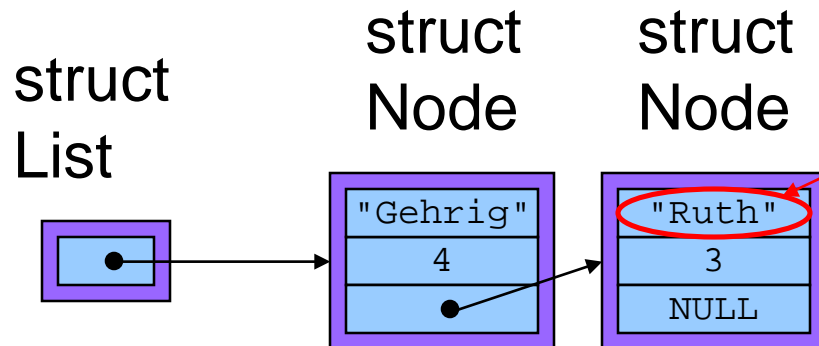
- (student name, grade)
  - (“john smith”, 84), (“jane doe”, 93), (“bill clinton”, 81)
- (baseball player, number)
  - (“Ruth”, 3), (“Gehrig”, 4), (“Mantle”, 7)
- (variable name, value)
  - (“maxLength”, 2000), (“i”, 7), (“j”, -10)



# Linked List Data Structure

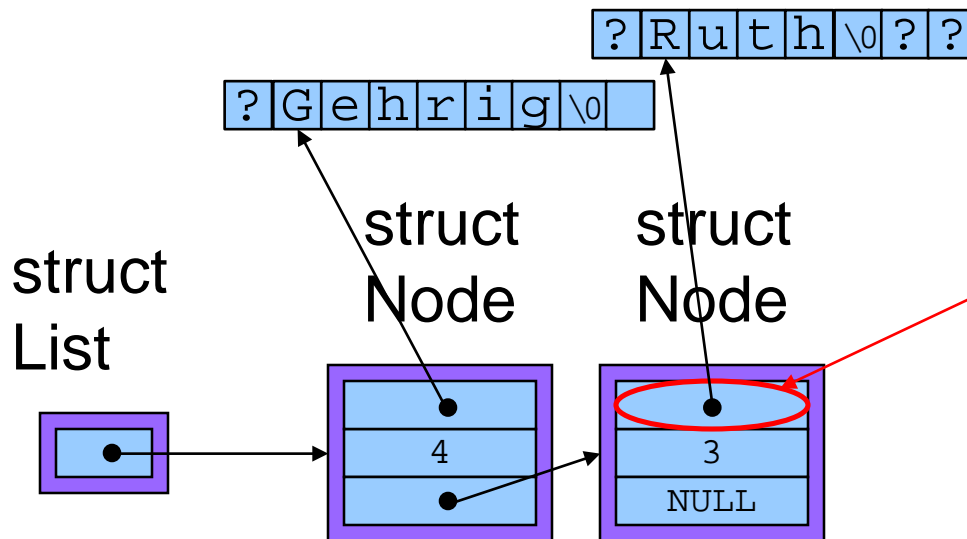
```
struct Node
{
  const char *key;
  int value;
  struct Node *next;
};

struct List
{
  struct Node *first;
};
```



Really this is the address at which "Ruth" resides

# Linked List Data Structure



Really this is the address at which "Ruth" resides



# Linked List Algorithms

## Create

- Allocate `List` structure; set `first` to `NULL`
- Performance:  $O(1) \Rightarrow$  fast

## Add (no check for duplicate key required)

- Insert new node containing key/value pair at front of list
- Performance:  $O(1) \Rightarrow$  fast

## Add (check for duplicate key required)

- Traverse list to check for node with duplicate key
- Insert new node containing key/value pair into list
- Performance:  $O(n) \Rightarrow$  slow



# Linked List Algorithms

## Search

- Traverse the list, looking for given key
- Stop when key found, or reach end
- Performance:  $O(n) \Rightarrow$  slow

## Free

- Free `Node` structures while traversing
- Free `List` structure
- Performance:  $O(n) \Rightarrow$  slow

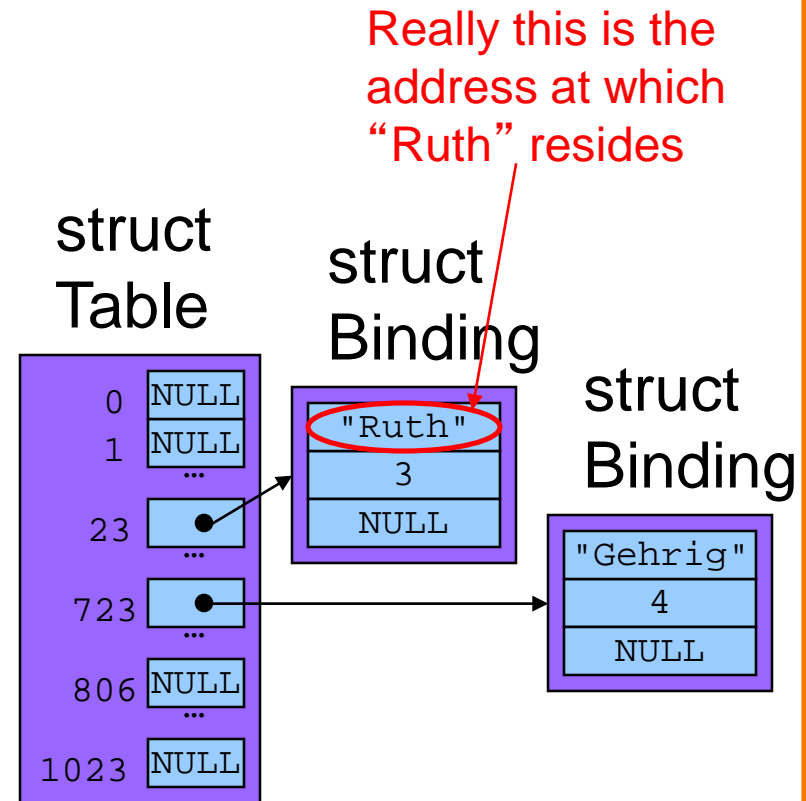
Would it be better to keep the nodes sorted by key?



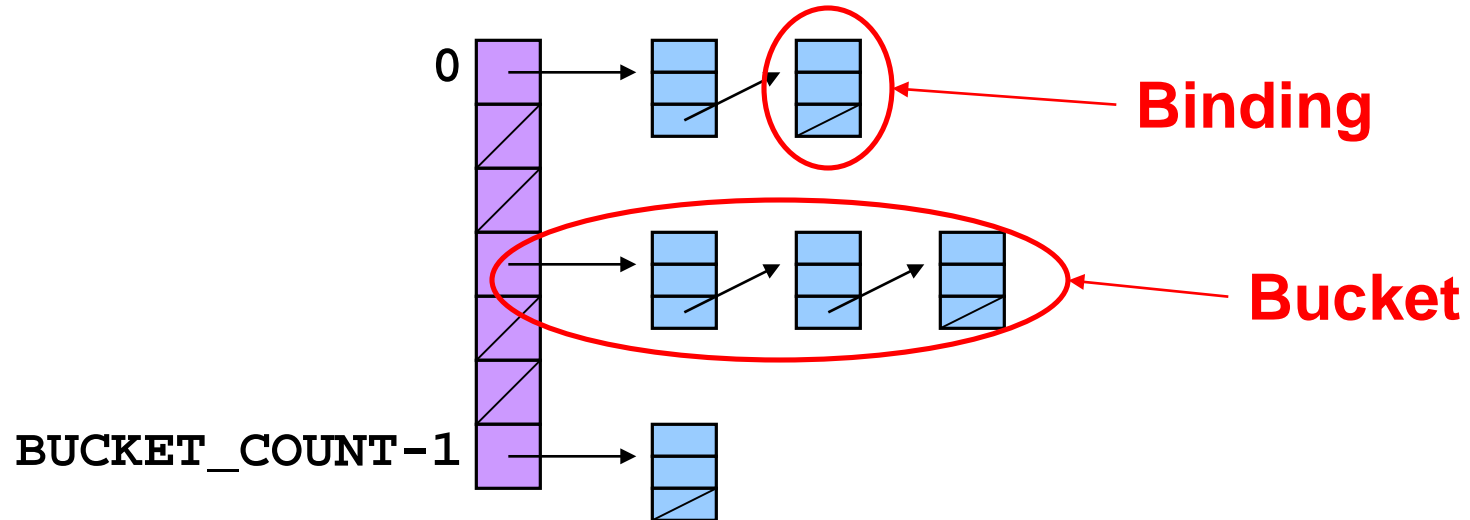
# Hash Table Data Structure

## Array of linked lists

```
enum {BUCKET_COUNT = 1024};  
  
struct Binding  
{  const char *key;  
   int value;  
   struct Binding *next;  
};  
  
struct Table  
{  struct Binding *buckets[BUCKET_COUNT];  
};
```



# Hash Table Data Structure



**Hash function** maps given key to an integer

Mod integer by `BUCKET_COUNT` to determine proper bucket





# Hash Table Example

Example: `BUCKET_COUNT = 7`

Add (if not already present) bindings with these keys:

- the, cat, in, the, hat

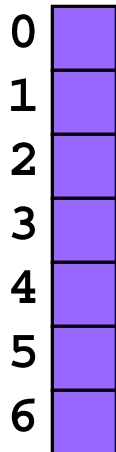


# Hash Table Example (cont.)

First key: “the”

- $\text{hash}(\text{“the”}) = 965156977$ ;  $965156977 \% 7 = 1$

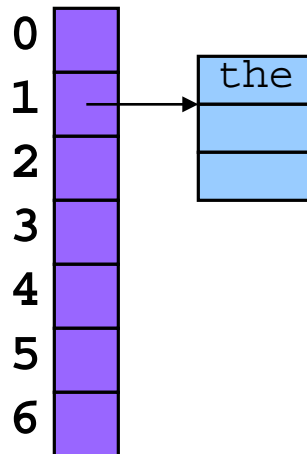
Search `buckets[1]` for binding with key “the”; not found



# Hash Table Example (cont.)



Add binding with key “the” and its value to `buckets[1]`



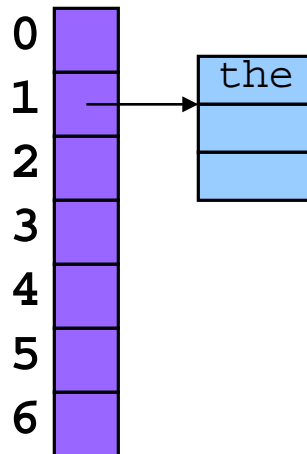


# Hash Table Example (cont.)

Second key: “cat”

- $\text{hash}(\text{“cat”}) = 3895848756; 3895848756 \% 7 = 2$

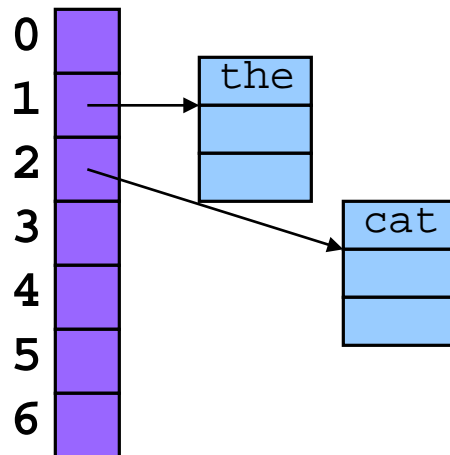
Search `buckets[2]` for binding with key “cat”; not found



# Hash Table Example (cont.)



Add binding with key “cat” and its value to `buckets[ 2 ]`



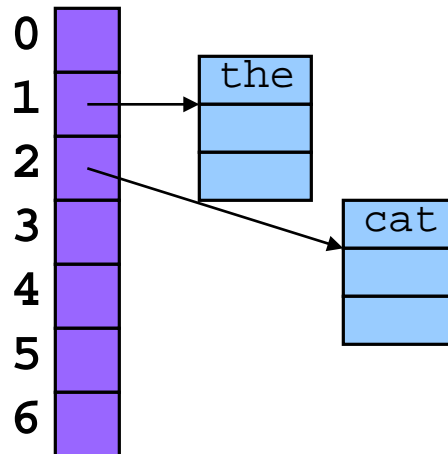


# Hash Table Example (cont.)

Third key: “in”

- $\text{hash}(\text{“in”}) = 6888005; 6888005 \% 7 = 5$

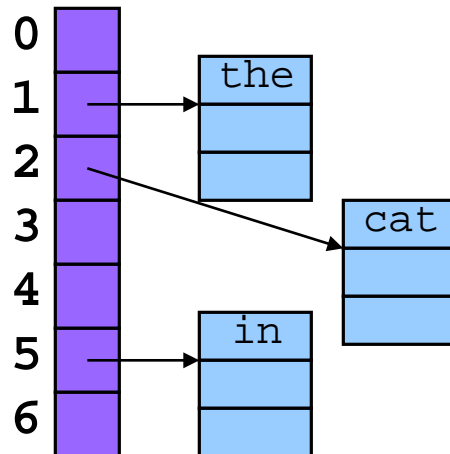
Search `buckets[5]` for binding with key “in”; not found



# Hash Table Example (cont.)



Add binding with key “in” and its value to `buckets[5]`





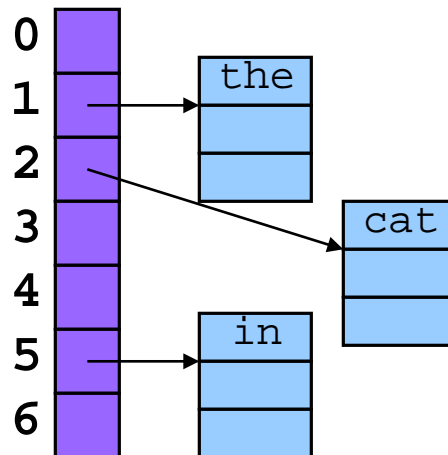
# Hash Table Example (cont.)

Fourth word: “the”

- $\text{hash}(\text{“the”}) = 965156977; 965156977 \% 7 = 1$

Search **buckets[1]** for binding with key “the”; found it!

- Don't change hash table





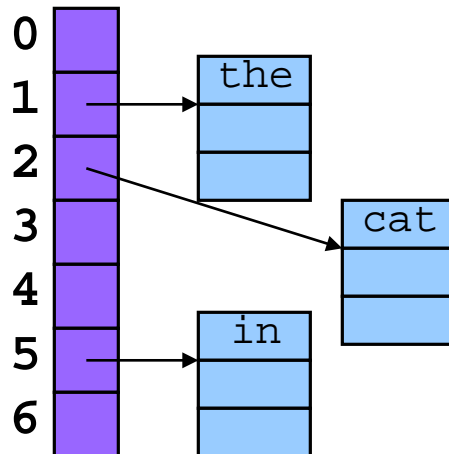


# Hash Table Example (cont.)

Fifth key: “hat”

- $\text{hash}(\text{“hat”}) = 865559739; 865559739 \% 7 = 2$

Search `buckets[2]` for binding with key “hat”; not found

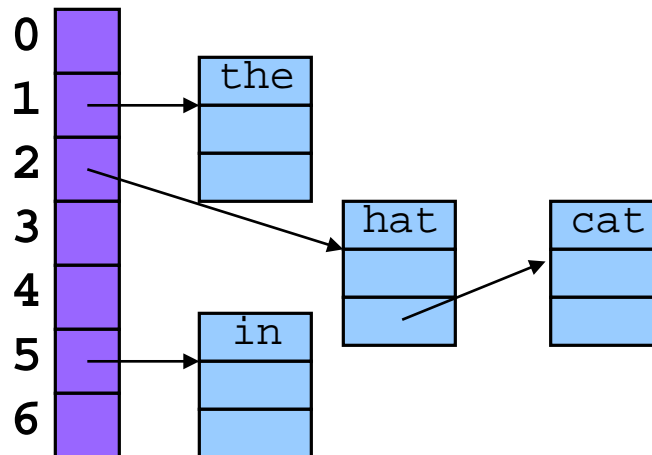




# Hash Table Example (cont.)

Add binding with key “hat” and its value to `buckets[ 2 ]`

- At front or back? Doesn't matter
- Inserting at the front is easier, so add at the front





# Hash Table Algorithms

## Create

- Allocate `Table` structure; set each bucket to `NULL`
- Performance:  $O(1) \Rightarrow$  fast

## Add

- Hash the given key
- Mod by `BUCKET_COUNT` to determine proper bucket
- Traverse proper bucket to make sure no duplicate key
- Insert new binding containing key/value pair into proper bucket
- Performance:  $O(1) \Rightarrow$  fast

Is the add performance always fast?



# Hash Table Algorithms

## Search

- Hash the given key
- Mod by **BUCKET\_COUNT** to determine proper bucket
- Traverse proper bucket, looking for binding with given key
- Stop when key found, or reach end
- Performance:  $O(1) \Rightarrow$  fast

## Free

- Traverse each bucket, freeing bindings
- Free **Table** structure
- Performance:  $O(n) \Rightarrow$  slow

Is the search performance always fast?



# How Many Buckets?

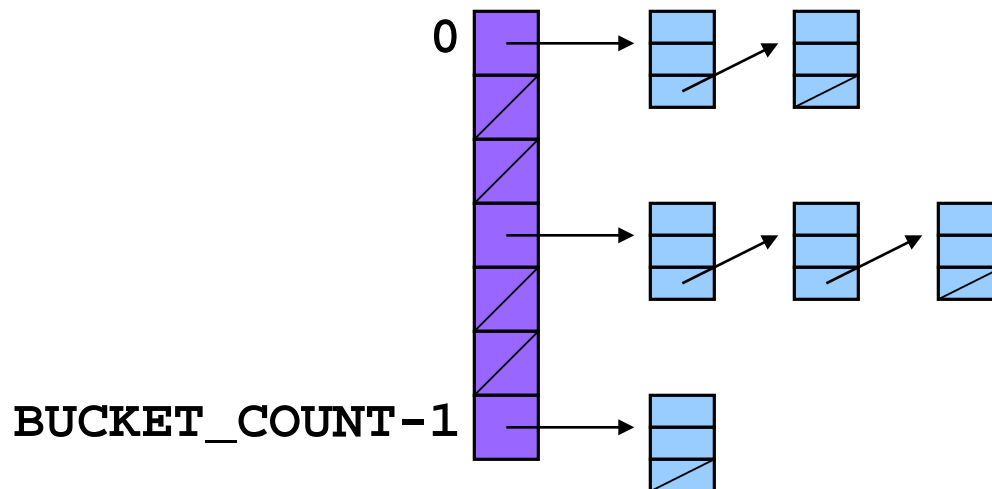
Many!

- Too few  $\Rightarrow$  large buckets  $\Rightarrow$  slow add, slow search

But not too many!

- Too many  $\Rightarrow$  memory is wasted

This is OK:

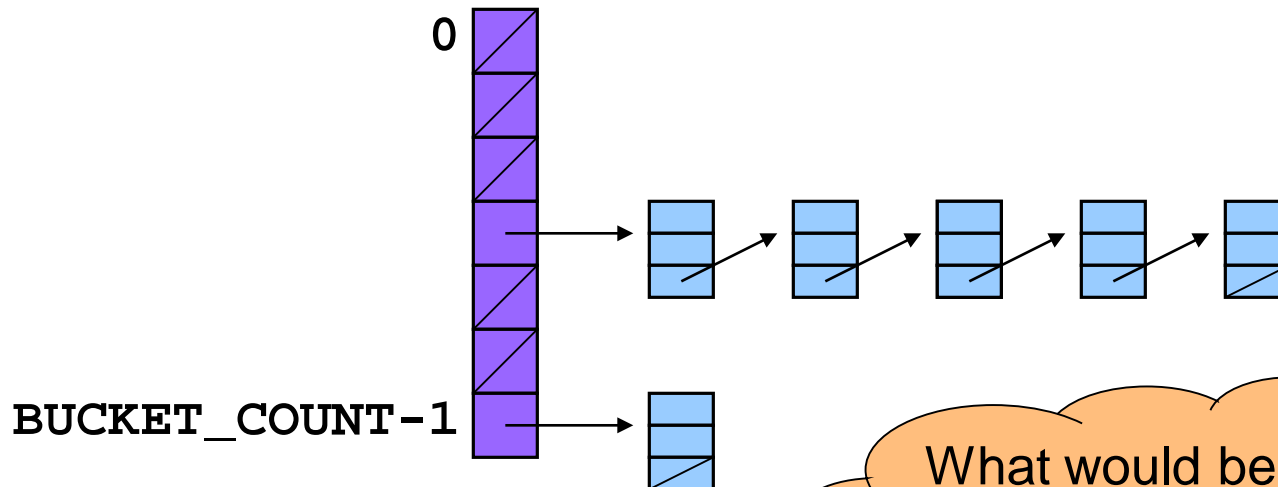


# What Hash Function?

Should distribute bindings across the buckets well

- Distribute bindings over the range  $0, 1, \dots, \text{BUCKET\_COUNT}-1$
- Distribute bindings *evenly* to avoid very long buckets

This is not so good:



What would be the worst possible hash function?



# How to Hash Strings?

Simple hash schemes don't distribute the keys evenly enough

- Number of characters, mod **BUCKET\_COUNT**
- Sum the numeric codes of all characters, mod **BUCKET\_COUNT**
- ...

A reasonably good hash function:

- Weighted sum of characters  $s_i$  in the string  $s$ 
  - $(\sum a^i s_i) \bmod \text{BUCKET\_COUNT}$
- Best if  $a$  and **BUCKET\_COUNT** are relatively prime
  - E.g.,  $a = 65599$ , **BUCKET\_COUNT** = 1024
- Even better if **BUCKET\_COUNT** is prime.

Why?

Footnote [A. Appel]: I originally designed this homework so that **BUCKET\_COUNT** is a prime number. In 2016 I wondered, "wouldn't it work just as well if  $a$  and **BUCKET\_COUNT** are just relatively prime? Measurements show no: using a prime number of buckets leads to more even distribution of bucket contents."



# How to Hash Strings?

Potentially expensive to compute  $\sum a^i s_i$

So let's do some algebra ("Horner's rule")

- (by example, for string  $s$  of length 5,  $a=65599$ ):

$$h = \sum 65599^i * s_i$$

$$h = 65599^0 * s_0 + 65599^1 * s_1 + 65599^2 * s_2 + 65599^3 * s_3 + 65599^4 * s_4$$

Direction of traversal of  $s$  doesn't matter, so...

$$h = 65599^0 * s_4 + 65599^1 * s_3 + 65599^2 * s_2 + 65599^3 * s_1 + 65599^4 * s_0$$

$$h = 65599^4 * s_0 + 65599^3 * s_1 + 65599^2 * s_2 + 65599^1 * s_3 + 65599^0 * s_4$$

$$h = (((((s_0) * 65599 + s_1) * 65599 + s_2) * 65599 + s_3) * 65599) + s_4$$



# How to Hash Strings?



Yielding this function

```
size_t hash(const char *s, size_t bucketCount)
{
    size_t i;
    size_t h = 0;
    for (i=0; s[i]!='\0'; i++)
        h = h * 65599 + (size_t)s[i];
    return h % bucketCount;
}
```

# How to Protect Keys?



Suppose `Table_add()` function contains this code:

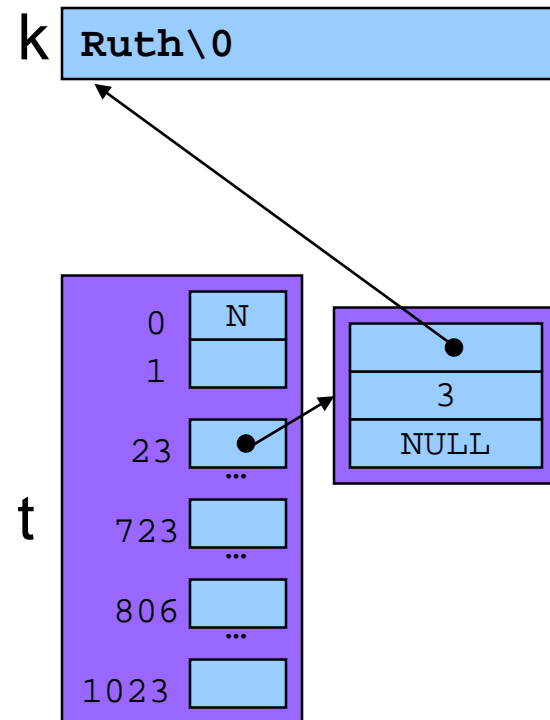
```
void Table_add(struct Table *t, const char *key, int value)
{
    ...
    struct Binding *p =
        (struct Binding*)malloc(sizeof(struct Binding));
    p->key = key;
    ...
}
```



# How to Protect Keys?

Problem: Consider this calling code:

```
struct Table *t;  
char k[100] = "Ruth";  
...  
Table_add(t, k, 3);
```



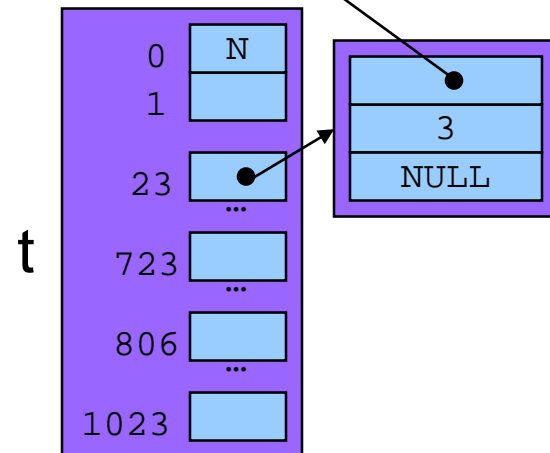


# How to Protect Keys?

Problem: Consider this calling code:

```
struct Table *t;  
char k[100] = "Ruth";  
...  
Table_add(t, k, 3);  
strcpy(k, "Gehrig");
```

k Gehrig\0



What happens if the client searches t for "Ruth"? For Gehrig?



# How to Protect Keys?

Solution: `Table_add()` saves a **defensive copy** of the given key

```
void Table_add(struct Table *t, const char *key, int value)
{
    ...
    struct Binding *p =
        (struct Binding*)malloc(sizeof(struct Binding));
    p->key = (const char*)malloc(strlen(key) + 1);
    strcpy((char*)p->key, key);
    ...
}
```

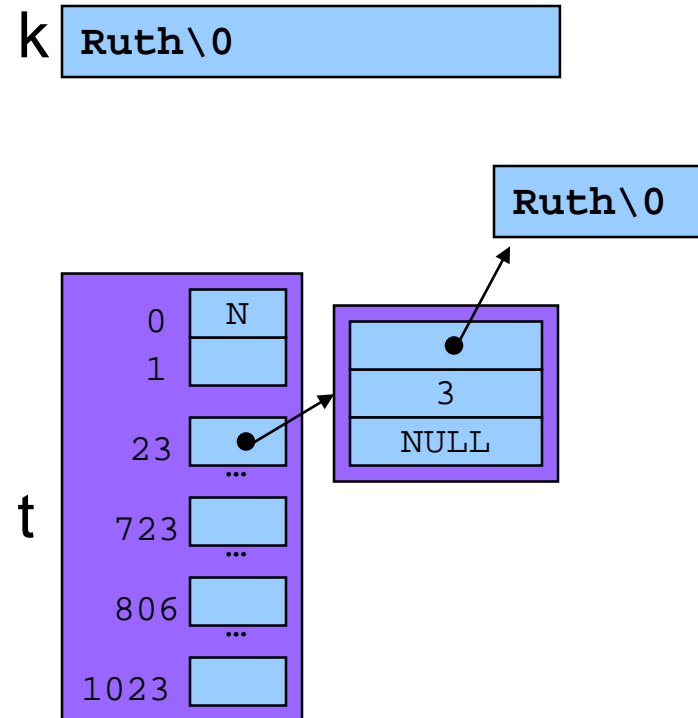
Why add 1?



# How to Protect Keys?

Now consider same calling code:

```
struct Table *t;  
char k[100] = "Ruth";  
...  
Table_add(t, k, 3);
```



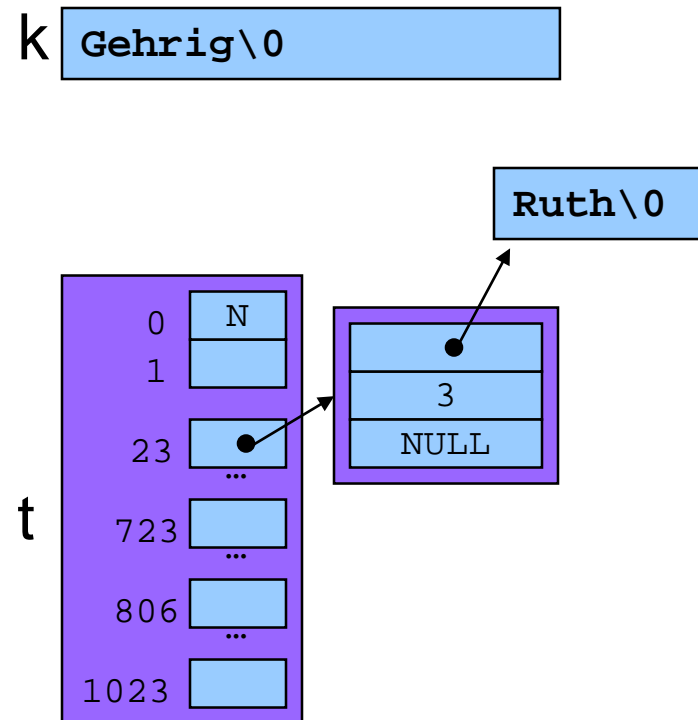


# How to Protect Keys?

Now consider same calling code:

```
struct Table *t;  
char k[100] = "Ruth";  
...  
Table_add(t, k, 3);  
strcpy(k, "Gehrig");
```

Hash table is  
not corrupted



# Who Owns the Keys?



Then the hash table **owns** its keys

- That is, the hash table owns the memory in which its keys reside
- `Hash_free( )` function must free the memory in which the key resides



# Summary



## Common data structures and associated algorithms

- Linked list
  - (Maybe) fast add
  - Slow search
- Hash table
  - (Potentially) fast add
  - (Potentially) fast search
  - Very common

## Hash table issues

- Hashing algorithms
- Defensive copies
- Key ownership