# Representing and Searching Sets of Strings

Antonio Carzaniga

Faculty of Informatics
Università della Svizzera italiana

April 27, 2016

- Radix search

- Ternary search tries

- Several very important applications

- Several very important applications

  E.g.,

  - dictionary (of words)
  - symbol table in a compiler
  - all kinds of key-based index
  - …

- Operations

- Operations

  - insert(Key)

  - delete(Key)

  - search(Key)

  - min()

  - max()

- Operations

- Operations

    - insert(Key)

    - search(Key)

- Operations

  - ▶ insert(Key)

  - ▶ search(Key)

- No *delete* operation

- Built once and searched many times

---

**BINARYSEARCH**($A$, $K$)

1   $first = 1$
2   $last = length(A)$
3   **while** $first \leq last$
4      $x = \lceil (first + last)/2 \rceil$
5      **if** $A[x] == K$
6         **return** TRUE
7      **elseif** $first == last$
8         **return** FALSE
9      **elseif** $A[x] > K$
10       $last = x - 1$
11     **else** $fisrt = x + 1$
12   **return** FALSE

---

**TREE-SEARCH**($T$, $K$)

1   $x = T.root$
2   **while** $x \neq$ NIL **and** $K \neq x.key$
3      **if** $K < x.key$
4         $x = x.left$
5      **else** $x = x.right$
6   **if** $x \neq$ NIL
7      **return** TRUE
8   **else return** FALSE

---

■ Complexity?

— *K **is a string!***

**BINARYSEARCH**($A$, $K$)

1  $first = 1$
2  $last = length(A)$
3  **while** $first \leq last$
4      $x = \lceil (first + last)/2 \rceil$
5      **if** $A[x] == K$
6          **return** TRUE
7      **elseif** $first == last$
8          **return** FALSE
9      **elseif** $A[x] > K$
10         $last = x - 1$
11     **else** $fisrt = x + 1$
12 **return** FALSE

**TREE-SEARCH**($T$, $K$)

1  $x = T.root$
2  **while** $x \neq$ NIL **and** $K \neq x.key$
3      **if** $K < x.key$
4          $x = x.left$
5      **else** $x = x.right$
6  **if** $x \neq$ NIL
7      **return** TRUE
8  **else return** FALSE

■ Complexity?

─── *K **is a string!*** ───

**BINARYSEARCH**($A$, $K$)

```
1   first = 1
2   last = length(A)
3   while first ≤ last
4       x = ⌈(first + last)/2⌉
5       if A[x] == K
6           return TRUE
7       elseif first == last
8           return FALSE
9       elseif A[x] > K
10          last = x − 1
11      else fisrt = x + 1
12  return FALSE
```

**TREE-SEARCH**($T$, $K$)

```
1   x = T.root
2   while x ≠ NIL and K ≠ x.key
3       if K < x.key
4           x = x.left
5       else x = x.right
6   if x ≠ NIL
7       return TRUE
8   else return FALSE
```

■ Complexity?

▶ we must account for the ***complexity of string comparisons***

- Assuming a string is an array of bytes, the condition $A[x]$ == $K$ (line 5 of **BINARYSEARCH**) becomes **STRINGEQUALS**($A[x]$, $K$)

# String Comparison

- Assuming a string is an array of bytes, the condition $A[x] == K$ (line 5 of **BINARYSEARCH**) becomes **STRINGEQUALS**$(A[x], K)$

**STRINGEQUALS**$(S_1, S_2)$
1   **if** $length(S_1) \neq length(S_2)$
2       **return** FALSE
3   **for** $i = 1$ **to** $length(S_1)$
4       **if** $S_1[i] \neq S_2[i]$
5          **return** FALSE
6   **return** TRUE

# String Comparison

- Assuming a string is an array of bytes, the condition $A[x] == K$ (line 5 of **BINARYSEARCH**) becomes **STRINGEQUALS**$(A[x], K)$

> **STRINGEQUALS**$(S_1, S_2)$
> 1  **if** $length(S_1) \neq length(S_2)$
> 2      **return** FALSE
> 3  **for** $i = 1$ **to** $length(S_1)$
> 4      **if** $S_1[i] \neq S_2[i]$
> 5          **return** FALSE
> 6  **return** TRUE

- The complexity of **STRINGEQUALS**$(S_1, S_2)$ is $O(m)$, where $m$ is the max string size

# String Comparison

- Assuming a string is an array of bytes, the condition $A[x] == K$ (line 5 of **BinarySearch**) becomes **StringEquals**$(A[x], K)$

```
StringEquals(S_1, S_2)
1  if length(S_1) ≠ length(S_2)
2      return FALSE
3  for i = 1 to length(S_1)
4      if S_1[i] ≠ S_2[i]
5          return FALSE
6  return TRUE
```

- The complexity of **StringEquals**$(S_1, S_2)$ is $O(m)$, where $m$ is the max string size

- So, the complexity of **BinarySearch**$(A, K)$ is $O(m \log n)$

**CHAINED-HASH-SEARCH**$(T, K)$

1  $L = T[h(K)]$
2  **return LIST-SEARCH**$(L, K)$

**HASH-SEARCH**$(T, K)$

1  **for** $i = 1$ **to** $length(T)$
2      $j = h(K, i)$
3      **if** $T[j] == K$
4          **return** TRUE
5      **if** $T[j] ==$ NIL
6          **return** FALSE
7  **return** FALSE

## What About a Hash Table

**CHAINED-HASH-SEARCH**$(T, K)$

1  $L = T[h(K)]$
2  **return LIST-SEARCH**$(L, K)$

**HASH-SEARCH**$(T, K)$

1  **for** $i = 1$ **to** $length(T)$
2      $j = h(K, i)$
3      **if** $T[j] == K$
4          **return** TRUE
5      **if** $T[j] ==$ NIL
6          **return** FALSE
7  **return** FALSE

- Complexity?

**CHAINED-HASH-SEARCH**$(T, K)$

1   $L = T[h(K)]$
2   **return LIST-SEARCH**$(L, K)$

**HASH-SEARCH**$(T, K)$

1   **for** $i = 1$ **to** $length(T)$
2       $j = h(K, i)$
3       **if** $T[j] == K$
4           **return** TRUE
5       **if** $T[j] ==$ NIL
6           **return** FALSE
7   **return** FALSE

■ Complexity?

▶ here, too, we must account for the string comparisons

**CHAINED-HASH-SEARCH**$(T, K)$

1  $L = T[\boxed{h(K)}]$
2  **return LIST-SEARCH**$(L, K)$

**HASH-SEARCH**$(T, K)$

1  **for** $i = 1$ **to** $length(T)$
2      $j = \boxed{h(K, i)}$
3      **if** $\boxed{T[j] == K}$
4          **return** TRUE
5      **if** $T[j] == $ NIL
6          **return** FALSE
7  **return** FALSE

- Complexity?
  - here, too, we must account for the string comparisons
  - and for the hash functions

- When we start **BINARYSEARCH**$(A, K)$
  - ▶ $A[x]$ is probably far away from $K$
  - ▶ so, **STRINGEQUALS**$(A[x], K)$ is likely to return quickly

- When we start **BINARYSEARCH**(*A*, *K*)

  ▶ *A*[*x*] is probably far away from *K*

  ▶ so, **STRINGEQUALS**(*A*[*x*], *K*) is likely to return quickly

- Later in **BINARYSEARCH**(*A*, *K*)

  ▶ *A*[*x*] gets closer and closer to *K*

  ▶ so, **STRINGEQUALS**(*A*[*x*], *K*) is likely to iterate for nearly *m* steps

- When we start **BinarySearch**$(A, K)$

  - $A[x]$ is probably far away from $K$

  - so, **StringEquals**$(A[x], K)$ is likely to return quickly

- Later in **BinarySearch**$(A, K)$

  - $A[x]$ gets closer and closer to $K$

  - so, **StringEquals**$(A[x], K)$ is likely to iterate for nearly $m$ steps

  - problem is, **StringEquals**$(A[x], K)$ is likely to go through the same prefix of $K$ many times

- So, since $m = \Theta(\log N)$, and **BinarySearch**$(A, K)$ uses $\Theta(\log N)$ comparisons each one running in $O(m)$:

$$T(N, m) = O(\log^2 N)$$

- Idea: a data structure where ***common prefixes*** are shared

- Idea: a data structure where *common prefixes* are shared

- Data structure useful for information re_trie_val (pronounced "try" to distinguish it from a tree…)

- Data structure useful for information re*trie*val (pronounced "try" to distinguish it from a tree…)

- Every node holds one character

- Data structure useful for information re*trie*val (pronounced "try" to distinguish it from a tree…)

- Every node holds one character

- Keys with the same prefix share a branch of the tree

- Data structure useful for information re*trie*val (pronounced "try" to distinguish it from a tree…)

- Every node holds one character

- Keys with the same prefix share a branch of the tree

- Keys are stored at (or just represented by) leaf nodes

- Data structure useful for information re*trie*val (pronounced "try" to distinguish it from a tree…)

- Every node holds one character

- Keys with the same prefix share a branch of the tree

- Keys are stored at (or just represented by) leaf nodes

- *Question:* how do we represent nodes and links?

- Data structure useful for information re*trie*val (pronounced "try" to distinguish it from a tree…)

- Every node holds one character

- Keys with the same prefix share a branch of the tree

- Keys are stored at (or just represented by) leaf nodes

- *Question:* how do we represent nodes and links?

  - one way would be to hold $|\Sigma|$ links

  - one for each character of the given alphabet $\Sigma$

1 = a
2 = b
3 = c

$\vdots$

26 = z

a

⋮

l

n

⋮

z

1 = a •
2 = b
3 = c

⋮

26 = z

$1 = a$
$2 = b$
$3 = c$
$26 = z$

a
l
n
z

a
t
z

**Radix Trie**

**Radix Trie**

**Radix Trie**

**Radix Trie**

**Radix Trie**

**Radix Trie**

**Radix Trie**

- Every element *x* has an array of links *x*. *links*
    - ▶ e.g., in "radix-256," an element represents a *byte* in a string (of bytes)

- Every element *x* has a *x*. *value* that is TRUE if that prefix corresponds to a string in the dictionary
    - ▶ this is to distinguish an entire word from a prefix

# Radix Search

- Every element *x* has an array of links *x*. *links*
  - e.g., in "radix-256," an element represents a *byte* in a string (of bytes)

- Every element *x* has a *x*. *value* that is TRUE if that prefix corresponds to a string in the dictionary
  - this is to distinguish an entire word from a prefix

```
RADIXSEARCH(Root, K)

1  n = Root
2  for i = 1 to length(K)
3      if n.links[K[i]] == NIL
4          return FALSE
5      else n = n.links[K[i]]
6  return n.value
```

- What is the complexity of Radix Search with a dictionary of *N* strings of up to *m* characters?

- What is the complexity of Radix Search with a dictionary of *N* strings of up to *m* characters?

$$T(N, m) = \Theta(m)$$

- What is the complexity of Radix Search with a dictionary of *N* strings of up to *m* characters?

$$T(N, m) = \Theta(m)$$

- What is the *space complexity?*

- What is the complexity of Radix Search with a dictionary of $N$ strings of up to $m$ characters?

$$T(N, m) = \Theta(m)$$

- What is the ***space complexity?***
    - ▶ first approximation:

$$S(N, m) = O(|\Sigma|mN)$$

- What is the complexity of Radix Search with a dictionary of *N* strings of up to *m* characters?

$$T(N, m) = \Theta(m)$$

- What is the ***space complexity?***
  - ▶ first approximation:

$$S(N, m) = O(|\Sigma|mN) \qquad S(N, m) = \Omega(|\Sigma| \log_{|\Sigma|} N)$$

■ What is the complexity of Radix Search with a dictionary of $N$ strings of up to $m$ characters?

$$T(N, m) = \Theta(m)$$

■ What is the *space complexity?*

▶ first approximation:

$$S(N, m) = O(|\Sigma|mN) \qquad S(N, m) = \Omega(|\Sigma| \log_{|\Sigma|} N)$$

▶ a better characterization (*Exercise:* figure this out!):

$$S(N, m) = \Theta\left(|\Sigma|\left[\frac{N-1}{|\Sigma|-1} + N\left(m - \frac{\log N}{\log|\Sigma|}\right)\right]\right)$$

- We do not represent a full array of links

- We do not represent a full array of links

- Instead, we represent a small binary "index" of the existing links

- We do not represent a full array of links

- Instead, we represent a small binary "index" of the existing links
  E.g., prefixes "xa", "xb", "xn", "xk", and "xs" might be represented as
  follows

- We do not represent a full array of links

- Instead, we represent a small binary "index" of the existing links
  E.g., prefixes "xa", "xb", "xn", "xk", and "xs" might be represented as follows

- *n.character* is the character at node *n*; i.e., the last character in the *prefix* represented by *n*

- *n. character* is the character at node *n*; i.e., the last character in the *prefix* represented by *n*

- *n. value* is the *value* to which *n* maps to; if the TST is a dictionary, then *n. value* is *true* iff the prefix represented by *n* is a key in the dictionary

- *n. character* is the character at node *n*; i.e., the last character in the *prefix* represented by *n*

- *n. value* is the *value* to which *n* maps to; if the TST is a dictionary, then *n. value* is *true* iff the prefix represented by *n* is a key in the dictionary

- A node *n* has three links
  - *n. lower* links to a node representing a "lower" character at the same position
  - *n. higher* links to a node representing a "higher" character at the same position
  - *n. equal* links to a node representing a character in the next position

"culture"

"culture"

"lugano"

"lugano"

"lunatic"

"lunatic"

"ciao"

"ciao"

"cappero"

"cappero"

"class"

"class"

"classic"

"classic"

"algorithm"

"algorithm"

"algo"

"algo"

```
TSTSEARCH(T, K)
 1  for i = 1 to |K|
 2      if i > 1
 3          T = T.equal
 4      while T ≠ NIL and K[i] ≠ T.character
 5          if K[i] < T.character
 6              T = T.lower
 7          else T = T.higher
 8      if T == NIL
 9          return FALSE
10  return n.value
```
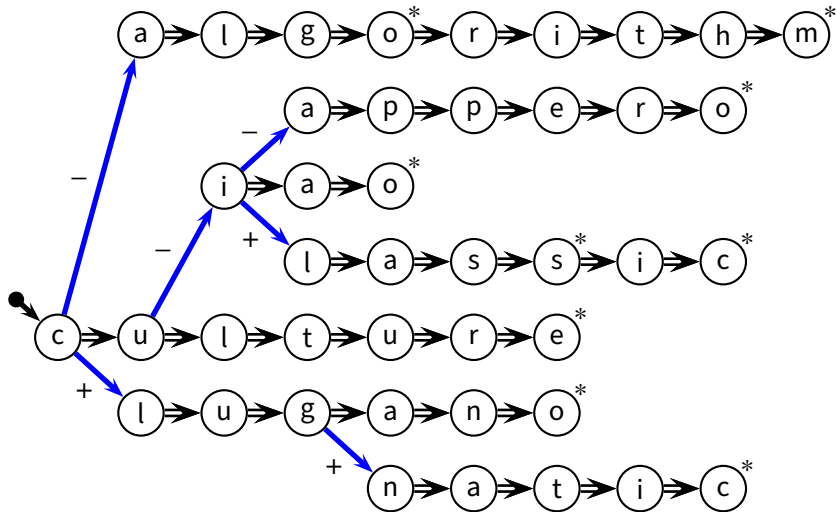
**TSTSEARCH**(*T*, *K*)
```
 1  for i = 1 to |K|
 2      if i > 1
 3          T = T.equal
 4      while T ≠ NIL and K[i] ≠ T.character
 5          if K[i] < T.character
 6              T = T.lower
 7          else T = T.higher
 8      if T == NIL
 9          return FALSE
10  return n.value
```

- Is it correct?

```
TSTSearch(T, K)
 1  for i = 1 to |K|
 2      if i > 1
 3          T = T.equal
 4      while T ≠ NIL and K[i] ≠ T.character
 5          if K[i] < T.character
 6              T = T.lower
 7          else T = T.higher
 8      if T == NIL
 9          return FALSE
10  return n.value
```

- Is it correct?    Not completely! (**Exercise:** fix it.)

- Complexity?

```
TSTSearch(T, K)
 1  for i = 1 to |K|
 2      if i > 1
 3          T = T.equal
 4      while T ≠ NIL and K[i] ≠ T.character
 5          if K[i] < T.character
 6              T = T.lower
 7          else T = T.higher
 8      if T == NIL
 9          return FALSE
10  return n.value
```

- Is it correct?    Not completely! (**Exercise:** fix it.)

- Complexity?    Non-trivial…

- Recursion starts with $root = $ **TSTINSERT**($root, K, 1$)

```
TSTINSERT(T, K, i)
 1  if T == NIL
 2      T = NEWNODE(K[i])
 3  if K[i] < T.character
 4      T.lower = TSTINSERT(T.lower, K, i)
 5  elseif K[i] > T.character
 6      T.higher = TSTINSERT(T.higher, K, i)
 7  elseif K[i] == T.character
 8      if i < |K|
 9          T.equal = TSTINSERT(T.equal, K, i + 1)
10      else T.value = TRUE
11  return T
```