# 6 Text Indexing –

## Searching whole genomes

*16 March 2021*

Sebastian Wild

# 6 Text Indexing

## 6.1 Motivation

# Text indexing

- *Text indexing* (also: *offline text search*):
    - case of string matching: find $P[0..m-1]$ in $T[0..n-1]$
    - but with *fixed* text $\rightsquigarrow$ preprocess $T$ (instead of $P$)
    - $\rightsquigarrow$ expect many queries $P$, answer them without looking at all of $T$
    - $\rightsquigarrow$ essentially a data structuring problem: "building an *index* of $T$"

        Latin: "one who points out"

- application areas
    - web search engines
    - online dictionaries
    - online encyclopedia
    - DNA/RNA data bases
    - . . . searching in any collection of text documents (that grows only moderately)

# Inverted indices

- original indices in books: list of (key) words $\mapsto$ page numbers where they occur

  same as "indexes"

- assumption: searches are only for **whole** (key) **words**

$\rightsquigarrow$ often reasonable for natural language text

# Inverted indices

same as "indexes"

- original indices in books: list of (key) words $\mapsto$ page numbers where they occur

- assumption: searches are only for **whole** (key) **words**

$\rightsquigarrow$ often reasonable for natural language text

**Inverted index:**

- collect all words in $T$

    - can be as simple as splitting $T$ at whitespace
    - ( actual implementations typically support *stemming* of words
      goes $\rightarrow$ go, cats $\rightarrow$ cat )

- store mapping from words to a list of occurrences $\rightsquigarrow$ *how?*

like a dictionary!    keys = words    ( BST
                                          but O(log n)
                                          time )

                     values = list of occurence,

2

# Clicker Question

Do you know what a *trie* is?

**A** A what? No!

**B** I have heard the term, but don't quite remember.

**C** I remember hearing about it in a module.
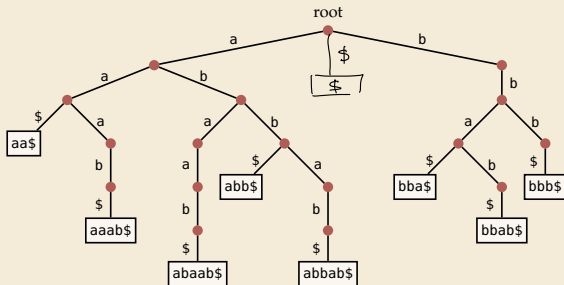
**D** Sure.

`sli.do/comp526`

**Click on "Polls" tab**

# Tries

- ▶ efficient dictionary data structure for strings

- ▶ name from re**trie**val, but pronounced "try"                    *≈ tree*

- ▶ tree based on symbol comparisons

- ▶ **Assumption:** stored strings are *prefix-free* (no string is a prefix of another)
  - ▶ strings of same length ✓
  - ▶ strings have "end-of-string" marker $ ✓          *some character ∉ Σ*

- ▶ **Example**:          $\Sigma = \{a, b\}$
  {aa\$, aaab\$, abaab\$, abb\$,
  abbab\$, bba\$, bbab\$, bbb\$ ~~}~~ *, \$}*

# Clicker Question

Suppose we have a trie that stores $n$ strings over $\Sigma = \{\texttt{A}, \ldots, \texttt{Z}\}$. Each stored string consists of $m$ characters.

We now search for a query string $Q$ with $|Q| = q$. $\quad (q \leq m)$

How many **nodes** in the trie are **visited** during this **query**?

**A** $\Theta(\log n)$        **F** $\Theta(\log m)$

**B** $\Theta(\log(nm))$        **G** $\Theta(q)$

**C** $\Theta(m \cdot \log n)$        **H** $\Theta(\log q)$

**D** $\Theta(m + \log n)$        **I** $\Theta(q \cdot \log n)$

**E** $\Theta(m)$        **J** $\Theta(q + \log n)$

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

Suppose we have a trie that stores $n$ strings over $\Sigma = \{\text{A}, \ldots, \text{Z}\}$. Each stored string consists of $m$ characters.

We now search for a query string $Q$ with $|Q| = q$.

How many **nodes** in the trie are **visited** during this **query**?

**A** $\Theta(\log n)$

**B** $\Theta(\log(nm))$

**C** $\Theta(m \cdot \log n)$

**D** $\Theta(m + \log n)$

**E** $\Theta(m)$

**F** $\Theta(\log m)$

**G** $\Theta(q)$ ✓

**H** $\Theta(\log q)$

**I** $\Theta(q \cdot \log n)$

**J** $\Theta(q + \log n)$

# Clicker Question

Suppose we have a trie that stores $n$ strings over $\Sigma = \{\texttt{A}, \ldots, \texttt{Z}\}$. Each stored string consists of $m$ characters.

How many **nodes** does the trie have **in total** _in the worst case_?

A $\quad \Theta(n)$

B $\quad \Theta(n + m)$

C $\quad \Theta(n \cdot m)$

D $\quad \Theta(n \log m)$

E $\quad \Theta(m)$

F $\quad \Theta(m \log n)$

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

Suppose we have a trie that stores $n$ strings over $\Sigma = \{\text{A}, \ldots, \text{Z}\}$. Each stored string consists of $m$ characters.

How many **nodes** does the trie have **in total** *in the worst case*?

**A** ~~$\Theta(n)$~~

**B** ~~$\Theta(n + m)$~~

**C** $\Theta(n \cdot m)$ ✓

**D** ~~$\Theta(n \log m)$~~

**E** ~~$\Theta(m)$~~
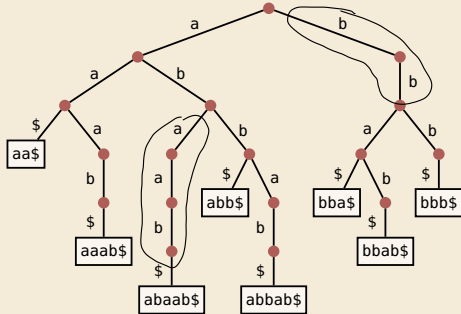
**F** ~~$\Theta(m \log n)$~~

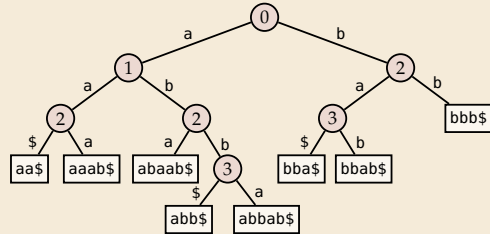`sli.do/comp526`

Click on "Polls" tab

# Compact tries

=1 child

► compress paths of unary nodes into single edge
► nodes store index of next character



**standard trie**

**compact trie**

⤳ searching slightly trickier, but same time complexity as in trie

not $O(n \cdot m)$

► all nodes ≥ 2 children  ⤳  #nodes ≤ #leaves = #strings  ⤳  linear space  $O(n)$

# Tries as inverted index

👍 simple

👍 fast lookup

👎 cannot handle more general queries:
- ▶ search part of a word
- ▶ search phrase (sequence of words)

# Tries as inverted index

👍 simple

👍 fast lookup

👎 cannot handle more general queries:
- search part of a word
- search phrase (sequence of words)

👎 what if the 'text' does not even have <u>words</u> to begin with?!
- biological sequences

  ```
  ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTGCTCTCCGGGGCCACGGCCACCGCTGCCCTGCCCCTGGAGGGTGGCCCCACCGGC
  CGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCCTGACTTTCCTCGCTTGGTGGTTTGAGTGGACCTCCCAGGC
  CAGTGCCGGGCCCCTCATAGGAGAGGAAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGCACCCCCCAGCAATCCGCGCGCCGGGACAGAA
  TGCCCTGCAGGAACTTCTTCTGGAAGACCTTCTCCTCCTGCAAATAAAACCTCACCCATGAATGCTCACGCAAGTTTAATTACAGACCTGAA
  ```

- binary streams

  ```
  000000101010011110101110000011111000111110111110011011010000111000100110111100000100011101010
  011011000011010110100000001000000001110101100000100001111010111011001000110010110111011011111
  110001010001011001010000001110101010011000000001101100001100111110000101 01010111011110000011
  101011100100101010101000001111101001100000001111001101010000000010010010000010110001100011011
  ```

↝ need new ideas